



OPEN

## The influence of human genetic variation on Epstein–Barr virus sequence diversity

Sina Rüeger<sup>1,2,40</sup>, Christian Hammer<sup>3,40</sup>, Alexis Loetscher<sup>2,4,40</sup>, Paul J. McLaren<sup>5,6</sup>, Dylan Lawless<sup>1,2</sup>, Olivier Naret<sup>1,2</sup>, Nina Khanna<sup>7</sup>, Enos Bernasconi<sup>8</sup>, Matthias Cavassini<sup>9</sup>, Huldrych F. Günthard<sup>10</sup>, Christian R. Kahlert<sup>11,12</sup>, Andri Rauch<sup>13</sup>, Daniel P. Depledge<sup>14</sup>, Sofia Morfopoulou<sup>14</sup>, Judith Breuer<sup>14</sup>, Evgeny Zdobnov<sup>2,4</sup>, Jacques Fellay<sup>1,2,15</sup>✉ & the Swiss HIV Cohort Study\*

Epstein–Barr virus (EBV) is one of the most common viruses latently infecting humans. Little is known about the impact of human genetic variation on the large inter-individual differences observed in response to EBV infection. To search for a potential imprint of host genomic variation on the EBV sequence, we jointly analyzed paired viral and human genomic data from 268 HIV-coinfected individuals with CD4 + T cell count < 200/mm<sup>3</sup> and elevated EBV viremia. We hypothesized that the reactivated virus circulating in these patients could carry sequence variants acquired during primary EBV infection, thereby providing a snapshot of early adaptation to the pressure exerted on EBV by the individual immune response. We searched for associations between host and pathogen genetic variants, taking into account human and EBV population structure. Our analyses revealed significant associations between human and EBV sequence variation. Three polymorphic regions in the human genome were found to be associated with EBV variation: one at the amino acid level (BRLF1:p.Lys316Glu); and two at the gene level (burden testing of rare variants in BALF5 and BBRF1). Our findings confirm that jointly analyzing host and pathogen genomes can identify sites of genomic interactions, which could help dissect pathogenic mechanisms and suggest new therapeutic avenues.

Human genetic variation plays a key role in determining individual responses after exposure to infectious agents. Even though susceptibility or resistance to a microbial challenge is the final result of dynamic interactions between host, pathogen and environment, human genetic polymorphisms have been shown to have an important, directly quantifiable impact on the outcome of various infections<sup>1,2</sup>.

Genome-wide association studies (GWAS) have proven powerful to identify genetic regions implicated in a wide range of complex traits in both health and disease<sup>3</sup>. In the field of infectious diseases, several clinical and laboratory phenotypes have been investigated, including, for example disease susceptibility<sup>4,5</sup>, clinical outcomes<sup>6</sup>, adaptive immunity<sup>7–9</sup> or drug response<sup>10</sup>. In chronically infected patients, however, the pathogen genome itself provides a promising complementary target to investigate the impact of host genomic diversity on infection. While one part of the variation observed in pathogen DNA or RNA sequence is present at the transmission event,

<sup>1</sup>School of Life Sciences, EPFL, Lausanne, Switzerland. <sup>2</sup>Swiss Institute of Bioinformatics, Lausanne, Switzerland. <sup>3</sup>Genentech Inc, 1 DNA Way, South San Francisco, CA, USA. <sup>4</sup>Department of Genetic Medicine and Development, University of Geneva Medical School, Geneva, Switzerland. <sup>5</sup>JC Wilt Infectious Diseases Research Centre, Public Health Agency of Canada, Winnipeg, MB, Canada. <sup>6</sup>Department of Medical Microbiology and Infectious Diseases, University of Manitoba, Winnipeg, MB, Canada. <sup>7</sup>Department of Biomedicine, University Hospital Basel, University of Basel, Basel, Switzerland. <sup>8</sup>Division of Infectious Diseases, Regional Hospital Lugano, Lugano, Switzerland. <sup>9</sup>Division of Infectious Diseases, University Hospital Lausanne, University of Lausanne, Lausanne, Switzerland. <sup>10</sup>Division of Infectious Diseases and Hospital Epidemiology, University Hospital Zurich, University of Zurich, Zurich, Switzerland. <sup>11</sup>Division of Infectious Diseases and Hospital Epidemiology, Cantonal Hospital St.Gallen, St.Gallen, Switzerland. <sup>12</sup>Childrens Hospital of Eastern Switzerland, St. Gallen, Switzerland. <sup>13</sup>Department of Infectious Diseases, Bern University Hospital, University of Bern, Bern, Switzerland. <sup>14</sup>Division of Infection and Immunity, University College London, London, UK. <sup>15</sup>Precision Medicine Unit, Lausanne University Hospital and University of Lausanne, Lausanne, Switzerland. <sup>40</sup>These authors contributed equally: Sina Rüeger, Christian Hammer and Alexis Loetscher. \*A list of authors and their affiliations appears at the end of the paper. ✉email: jacques.fellay@epfl.ch

another fraction is acquired during the course of an infection, resulting at least partially from selective pressure exerted by the host response on the infectious agent. The phenomenon of within-host evolution has been extensively investigated for both viruses<sup>11–13</sup> and bacteria<sup>11,14</sup>. Pathogen genomic variation can thus be considered an intermediate phenotype that is detectable as a footprint of within-host evolution. This can serve as a basis for a joint association analyses of host and pathogen genome variation, which we called genome-to-genome (G2G) analysis<sup>15</sup>, a more powerful approach than using a clinical outcome alone. A global description of the adaptive forces acting on a pathogen genome during natural infection holds the potential to identify novel therapeutic and diagnostic targets and could inform vaccine design efforts<sup>16</sup>.

A G2G analysis for the quickly evolving human immunodeficiency virus (HIV) identified strong associations of single nucleotide polymorphisms (SNPs) in the HLA class I region with multiple amino acid variants across the viral genome<sup>15</sup>. More recent work showed an impact of variation in the HLA class II and interferon lambda 4 (IFNL4) loci on hepatitis C virus (HCV) sequence diversity<sup>17–19</sup>. While the rate of evolutionary change in RNA viruses is higher than in DNA viruses<sup>20</sup>, the latter also present considerable amounts of inter- and intra-host variation. Among herpesviruses, it has been shown that human cytomegalovirus (HCMV) has higher genomic variability than other DNA viruses<sup>21</sup>. Recent genome sequencing efforts demonstrated that the same holds true for Epstein–Barr virus (EBV)<sup>22,23</sup>.

EBV is a widespread human pathogen that causes infectious mononucleosis in about 10% of individuals during primary infection. EBV infection occurs most often early in life, with about 30% of children being seropositive by age 5, 50% by age 10 and up to 80% by age 18<sup>24</sup>. This human infecting herpesvirus has also been associated with post-transplant lymphoproliferative disease<sup>25</sup> and could play a role in some autoimmune diseases<sup>26–29</sup>. In addition, EBV has oncogenic properties and is implicated in the pathogenesis of multiple cancer types, predominantly Burkitt's lymphoma, Hodgkin's and non-Hodgkin's lymphoma, nasopharyngeal carcinoma and gastric carcinoma<sup>30,31</sup>. More than 5% of the 2 million infection-associated new cancer cases in 2008 could be attributed to EBV<sup>32</sup>; it was also estimated to have caused 1.8% of cancer deaths in 2010, i.e. more than 140,000 cases<sup>33</sup>.

The EBV genome is approximately 170 Kbp long and encodes at least 80 proteins, not all of which have been definitively identified or characterized. After primary infection, the EBV genome persists in B cells as multicopy episomes that replicate once per cell cycle. In this latent mode, only a small subset of viral genes is expressed. Latent EBV can then reactivate to a lytic cycle, which involves higher gene expression and genome amplification for packaging into new infectious viral particles<sup>34</sup>.

A small number of host and viral genomic analyses of EBV infection have been recently published, demonstrating that human and pathogen genetic diversity plays a role in disease outcome. A study in 270 EBV isolates from southern China identified two non-synonymous EBV variants within the *BALF2* gene that were strongly associated with the risk of nasopharyngeal carcinoma<sup>35</sup>. Another group investigated the co-evolution of worldwide EBV strains<sup>36</sup> and found extensive linkage disequilibrium (LD) throughout EBV genomes. Furthermore, they observed that genes in strong LD were enriched in immunogenic genes, suggesting adaptive immune selection and epistasis. In a pediatric study of 58 Endemic Burkitt lymphoma cases and 40 healthy controls, an EBV genome GWAS identified 6 associated variants in the genes *EBNA1*, *EBNA2*, *BcLF1*, and *BARF1*<sup>37</sup>. On the population genetic side, a study of > 150 EBV genomes with known geographical origin revealed considerable variation in allele frequencies of EBV sub-populations<sup>38</sup>. Finally, the narrow-sense heritability of the humoral immune response against EBV was estimated to be 0.28<sup>9,39</sup>.

Here, we present the first global analysis of paired human and EBV genomes. We studied full EBV genomes together with their respective host genomic variation in a cohort of 268 immunocompromised, HIV-coinfected patients. We chose untreated HIV-coinfected patients because EBV reactivation leading to viremia is more prevalent in immunosuppressed individuals than in an average population. Our analysis reveals three novel host genomic loci that are associated with variation in EBV amino acids or genes.

## Materials and methods

**Study participants, sample preparation.** The Swiss HIV Cohort Study (SHCS) is a nationwide, prospective cohort study of HIV-infected patients that enrolled > 20,000 individuals since its establishment in 1988 and prospectively followed them at 6-month intervals<sup>40</sup>. For this project, SHCS participants were identified based on written consent for human genetic testing and availability of a peripheral blood mononuclear cell (PBMC) sample at time of advanced immunosuppression (i.e., with CD4 + T cell count below 200/mm<sup>3</sup>) in the absence of antiretroviral treatment.

We obtained demographic and clinical information from the SHCS database. These included sex, age, longitudinal HIV viral load results (number of RNA copies per ml of plasma), longitudinal CD4 + T cell counts (number of cells per mm<sup>3</sup> of blood), and history of opportunistic infections.

The SHCS has been approved by the Ethics Committees of all participating institutions (Ethikkommission Nordwest- und Zentralschweiz, EKNZ; Kantonale Ethikkommission Bern; Ethikkommission Ostschweiz, EKOS; Ethikkommission Zürich; Commission cantonale d'éthique de la recherche sur l'être humain, Genève, CCER; Commission cantonale d'éthique de la recherche sur l'être humain, Vaud, CER-VD; Comitato etico cantonale Ticino). Each study participant provided written informed consent for genetic testing, and all research was performed in accordance with relevant guidelines and regulations.

**EBV genome quantification, enrichment and sequencing.** DNA was extracted from PBMCs using the MagNA Pure 96 DNA and the Viral NA Small Volume Kit (Roche, Basel, Switzerland). Cellular EBV load was then determined using quantitative real-time PCR. Samples that yielded > 100 viral copies/ul were selected for EBV genome sequencing.

We used the previously described enrichment procedure to increase the relative abundance of EBV compared to host DNA<sup>41</sup>. Shortly, baits covering the EBV type 1 and 2 reference genomes were used to selectively capture viral DNA according to the SureSelect Illumina paired-end sequencing library protocol. Samples were then multiplexed and sequenced on an Illumina NextSeq sequencer<sup>41</sup>.

**EBV sequence analyses.** We chose a reference-based approach to call variants in the pathogen data. Since EBNA-2 and EBNA-3s are highly variable between EBV-1 and EBV-2 strains, we suspected that reads sequenced from these genes would map only to their corresponding type. In an attempt to attenuate the reference-bias this could cause, we constructed two references, one with the whole genome of the EBV-1 strain B95-8 (accession NC\_007605) and EBNA-2 and EBNA-3s sequences from EBV-2 strain AG876 (accession NC\_009334) and another one with the whole genome of AG876 with the EBNA-2 and EBNA-3s sequences of B95-8.

The read libraries were processed through Trimmomatic<sup>42</sup> to remove remnant PCR tags, TagDust<sup>43</sup> to eliminate low complexity reads and CD-HIT<sup>44</sup> to filter out duplicate reads. The remaining sequence reads were aligned to the constructs described in the previous paragraph. Following GATK best practices<sup>45,46</sup>, we mapped the read libraries using BWA mem<sup>47</sup>. We cleaned the regions around InDels using GATK v3.8's IndelRealigner<sup>48</sup>. As a last pre-processing step, we applied bwa-postalt.js, a BWA script that adjusts mapping quality score in function of alignments on ALT haplotypes.

Because patients can be infected by multiple EBV strains<sup>49</sup>, we used BWA's ALT-aware ability. In short, reads mapping to an ALT contig were always marked as supplementary alignment, regardless of mapping quality, unless they did not map to the primary assembly. This makes it easy to find unambiguously mapped reads, which we used as markers to quantify type 1 and type 2 EBV reads in all samples.

$$r1 = \frac{\#T1}{(\#T1 + \#T2) \cdot L1}$$

$$r2 = \frac{\#T2}{(\#T1 + \#T2) \cdot L2}$$

$$r = r1 - r2$$

where  $\#T1$  and  $\#T2$  are the unambiguous read counts against type 1 and 2 haplotypes, respectively,  $L1$  and  $L2$  are the length of type 1 and 2 haplotypes, respectively, and  $r1$  and  $r2$  are the type 1 and 2 ratios, respectively. The score  $r$  is the relative abundance between type 1 and type 2.

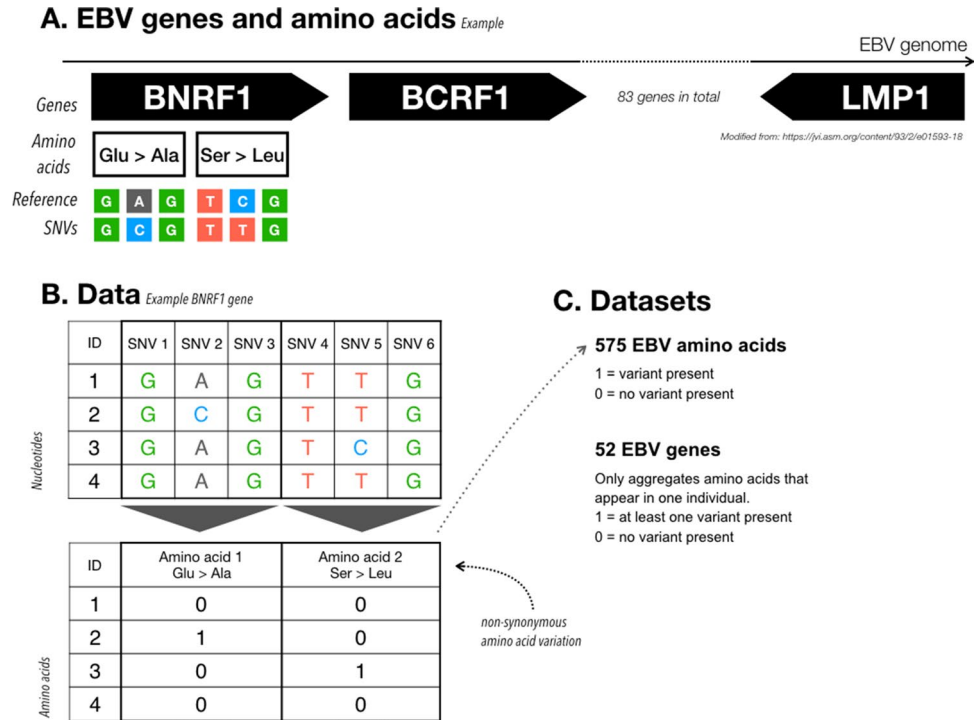
**Definition of EBV amino acid variants.** Since no gold standard variant set exists for EBV nor any closely related viral species, variant calling was performed using three different variant callers (GATK haplotypcaller, SNVer<sup>50</sup> and VarScan2<sup>51</sup>) and by selecting as *bona fide* variant set the intersection of the three. The identified EBV variants were annotated using snpEff<sup>52</sup>. Nucleotide variants were transformed into binary amino acid matrices using in-house Python scripts. The whole pipeline is written in Snakemake<sup>53</sup> and Python<sup>54</sup>.

This approach was benchmarked using synthetic libraries generated from B95-8 and AG876 using ART Illumina and RNFTools 0.3.1.3<sup>55</sup>, at a range of coverage between 10 and 250X and 5 different admixture conditions, 100% B95-8 or AG876, 75–25% and 50–50%. Assessing the true number of variants between EBV-1 and EBV-2 strains is not trivial because of the high variability in EBNA-2 and EBNA-3s regions. Therefore, we rated the variant callers and the consensus of the three mentioned callers on self-consistency. The performances of the runs were measured using the ratio of the variant counts to the size of the union of all variants called by a specific tested tool.

By using EBV type 2 as a reference, we focused on two types of variation in EBV strains: (1) single amino acid variants; and (2) burden of very rare amino acid variants (present in only 1 sample) in each viral gene (Fig. 1). We call these datasets *EBV amino acids* and *EBV genes*, respectively. Both datasets contain binary values, with a value of 1 standing for "variant present" and 0 for "no variant present". Positions with a coverage of less than 6× were set to "missing" and samples with more than 80% missing positions were excluded entirely. The positions covered by less than 6 reads were considered missing and imputed using the imputePCA function implemented in the missMDA R-package<sup>56</sup>. In total, we obtained 4392 amino acid variants and 83 gene variants. However, to limit the risk of model overfitting and because of low statistical power due to sample size we only included in the downstream association analyses the 575 *EBV amino acids* with an amino acid frequency of more than 10% and 52 *EBV genes*.

**Human genotyping and imputation.** A subset of 84 participants had been genotyped in the context of previous studies on several platforms. For the remaining 196 samples, human genomic DNA was isolated from PBMCs with the QIAAsymphony DSP DNA Kit (Qiagen, Hilden, Germany), and genotyped using Illumina OmniExpress (v1.1) BeadChip arrays.

Genotype imputation was performed on the Sanger imputation server independently for all genotyping platforms, using EAGLE2<sup>57</sup> for pre-phasing and PBWT<sup>58</sup> with the 1000 Genomes Phase 3 reference panel<sup>59</sup>. Low-quality imputed variants were excluded based on imputation INFO score (<0.8). All datasets were merged, only keeping markers that were genotyped or imputed for all genotyping platforms. SNPs were excluded on the basis of per-individual missingness (>3%), genotype missingness (>1%), marked deviation from Hardy–Weinberg equilibrium ( $p < 1 \times 10^{-6}$ ) and minor allele frequency <5% (Table 1). All quality control procedures were performed using PLINK 2.0<sup>60</sup>.



**Figure 1.** Illustration of EBV sequence variation. (A) The EBV genome is about 170 Kbp long and contains 83 genes, for a total of 4392 amino acid residues. As an example, we focus on the *BNRF1* gene and on two amino acid changes: Glu → Ala and Ser → Leu. We know for each sample the genomic variants across the whole genome, as illustrated with the colored nucleotides. Using the nucleotide information and a reference genome we can compute the amino acid changes. (B) We compare each individual (ID) to reference data and encode an amino acid as 1 if that individual has a non-synonymous change, and a 0 if not. This process returns us a matrix containing binary values, with individuals as row, and amino acids as columns. In our example, individual 2 has an amino acid change Glu → Ala and individual 3 an amino acid change Ser → Leu. (C) To transform the data into outcomes for the G2G analysis we can use the amino acid matrix as it is (*EBV amino acids* dataset) or remove all amino acid columns that appear in more than 1 individual and then pool amino acids per gene (1 = variant present, *EBV genes* dataset).

Dataset	Variable	Counts	Mean	Median	SD	Min	Max
Pathogen genome (83 Rare EBV gene variation)	Variant frequency		0.067	0.049	0.057	0	0.37
Pathogen genome (575 EBV amino acids)	Variant frequency		0.26	0.24	0.12	0.093	0.5
Covariates	Sex	Male: 206, female: 62					
Covariates	AGE		42.02	40.79	10.98	20.25	77.52
Covariates	PC1		-8.95	-7.12	32.53	-69.95	76.11
Covariates	PC2		7.39	8.4	30.41	-63.83	61.8
Covariates	PC3		4.45	3.19	21.41	-46.44	51.11
Covariates	PC4		-3.91	-6.69	20.72	-36.49	52.11
Covariates	PC5		-5.09	-8.39	18.76	-51.7	56.02
Covariates	PC6		-3.16	-3.21	18.65	-55.69	35.68
Covariates	EBV type		0.54	0.95	0.67	-1	1

**Table 1.** Summary of pathogen variants, host SNPs and covariates for 268 individuals. For each covariate, we indicate the number of individuals measured, and distribution (mean, median, standard deviation, minimum, maximum for quantitative, frequency for sex). For aggregated EBV genes the frequency is shown, for host SNPs the MAF distribution is presented.

**Association analyses.** We used the mixed model association implementation for binary and continuous outcomes in GCTA (v1.92)<sup>61,62</sup> to search for potential associations between human SNPs and EBV variants. The model can be expressed with the following equation:

$$y_k = \alpha X + \beta^{(kl)} g_l + \eta + \varepsilon$$

where the outcome  $y$  is a binary vector indicating whether an EBV variant is present (1) or not (0);  $X$  is a matrix that contains all covariates,  $\alpha$  represents all fixed effects of all covariates (including an intercept term),  $g$  is the SNP genotype vector with coded additive allele dosages 0, 1 or 2,  $\beta$  is the (fixed) effect of the SNP to be tested for association,  $\eta$  is the polygenic (random) effect and  $\varepsilon$  the error term. This mixed model was estimated for each EBV variant ( $k$ ) and SNP ( $l$ ), and integrated over all  $L$  SNPs and  $K$  EBV variants. To estimate  $\eta$ , the host genetic relationship matrix (GRM) was calculated from QC preprocessed genotype data using GCTA<sup>61</sup>.

The use of a mixed effects association model allows to account for population stratification of the host genome. To control for population stratification among EBV genomes, we included the first six principal components (PCs) of EBV genetic variation to the covariate matrix  $X$ <sup>63</sup>. Other covariates were sex, age and EBV type. PCs were calculated from EBV amino acid variants using the `convexLogisticPCA` function from the R package `logisticPCA`<sup>64</sup> in R<sup>65</sup>. As data preparation for PC computation, we removed variants with less than 5% or more than 95% frequency. Missing amino acid values were imputed with the `imputePCA` function from the R package `missMDA`<sup>66</sup>.

Significance was assessed using the usual genome-wide significance threshold in European populations of  $5 \times 10^{-8}$  and dividing it by the effective number of GWASs performed<sup>66</sup>. We used FINEMAP<sup>67</sup> to determine the most likely causal SNP(s) in a 2-Mb-wide window around each significant SNP. FINEMAP requires GWAS summary statistics and LD estimations as input. To estimate LD between SNPs, we used LDstore<sup>68</sup>. We performed eQTL lookups for host SNPs in eQTLGen<sup>69</sup>, EUGENE<sup>70</sup> and GTEx<sup>71</sup>.

Unless otherwise specified, all data preparation and analyses were performed using R<sup>65</sup>.

## Results

**Study participants and human genetic data.** PBMC samples from 778 SHCS participants were screened for the presence of cellular EBV DNA using RT-PCR. A total of 290 of them were identified as viremic for EBV (>2000 copies). We obtained good quality human genotyping and EBV sequencing data for 268 of them, which were included in the association analyses. The study cohort comprised 206 male and 62 female individuals, between the ages of 20 and 78 (median 40) (Table 1 and Supplementary Figure S1).

We applied standard GWAS quality control (QC) procedures that yielded information for 4'291'179 SNPs (Table 1 and Supplementary Figure S2, which shows the distribution of the minor allele frequency spectrum after QC).

**EBV genomic diversity and variant calling.** Genome coverage was very uneven between the samples. Mean depth varied from less than  $6 \times$  for 14 samples, up to more than  $500 \times$  in 5 others. We also observed fluctuation in coverage above  $6 \times$ , which we used to exclude 12 samples in which less than 20% of the EBV genome was sufficiently covered (Supplementary Figure S7a). In addition, the coverage in the first sequencing batch was not uniform.

We estimated the clonality of EBV strain in each sample by taking advantage of the high divergence between EBNA<sub>s</sub> T1 and T2 haplotypes. Among the 282 sequenced samples, 57.1% were predominantly (9:1) infected by T1 EBV, while 5.7% were mostly infected by T2 EBV (Supplementary Figure S6). The remaining 37.2% were infected by multiple strains or by recombinant viruses. This approach does not allow to stratify further than the EBNA types.

The variant calling pipeline was adapted to output variants by minimizing the impact of the admixture ratio and of the low coverage observed in the SHCS samples. Variants were called against EBV-2, as EBV-2 was able to call more variants than EBV-1 (Supplementary Figure S7d). The benchmark experiments against AG876 (EBV-2) yielded a total of 961 different variants. The most conservative was SNVer (783 variants), while the most sensitive was BCFtools 1.10.2-9<sup>72</sup> (930 variants). The variant callers can be prone to artifacts<sup>73</sup>, which was specifically observed in SNVer (Supplementary Figure S7c) in these datasets. To reduce the probability of calling artifacts, we chose to use the *bona fide* intersection of GATK HC, SNVer and VarScan2. This approach is likely to be impacted by low coverage. The recall is stable at around 95% at 25X coverage upwards and reasonable (10%) at 20X (Supplementary Figure S7c). Hence, low coverage has an impact, specifically, half potential variants called, on only 15% of the SHCS sample. However, this approach is very conservative, since it outputs only 88%, 85% and 79% of the variants called by SNVer, GATK HC and VarScan2, respectively.

On average, around 800 amino acid variants were called for each sample, with slight differences correlating with the clonality of the samples and the coverage above  $6 \times$  (Supplementary Figure S7d). The variant counts against the AG876 construct (EBV-2) were generally higher in mixed infections and EBV-1 strains (Supplementary Figure S7d A). The variant counts were generally lower in the samples included in the first sequencing batch, which is likely due to the fluctuating coverage. However, overall, the number of variants was found to be comparable across the samples, ranging from 400 to 1500 for the 77% samples with a  $6 \times$  coverage above 80% (Supplementary Figure S7d). Under 80% coverage, the variant counts hardly exceed 500 but rarely drops under 200 either. It is therefore likely that we missed variants using our approach. The positions covered by less than 6 reads were considered missing and imputed afterwards using the `imputePCA` function implemented in the `missMDA` R-package.

We analyzed EBV variation using two approaches: single marker analysis of *EBV amino acids*, to investigate common viral variation, and burden testing of very rare amino acid variants in *EBV genes* (Table 1, Supplementary Figure S3). Applying logistic principal component analysis of viral genomic structure showed a single main cluster (Supplementary Figure S4).

Locus	EBV dataset	EBV outcome	SNP	Chr	OR**	p	Finemapped SNP	Finemapped probability	Effect allele	EAF	n	Gene	Consequence type
2	Gene (binary, variants < 1 sample)	BALF5	rs2950922*	8	1.30739	4.2E-12	TRUE	0.13317	G	0.10821	268	UNC5D***	Intron_variant
3	Gene (binary, variants < 1 sample)	BBRF1	rs62124869*	2	1.29103	4.3E-11	TRUE	0.71847	C	0.0541	268	LINC01830	Non_coding_transcript_variant
1	Amino acid (binary)	BRLF1:p.Lys316Glu	rs7808072*	7	1.41346	6.8E-11	FALSE		T	0.06762	244		
1	Amino acid (binary)	BRLF1:p.Lys316Glu	rs6466720	7	1.42464	7.7E-11	TRUE	0.04783	G	0.06967	244		

**Table 2.** Summary of G2G analysis results. Top SNP and/or fine-mapped SNP per locus represented with: EBV dataset, EBV outcome, chromosome, SNP identifier, odds ratio, *p* value, whether this SNP is a top SNP or a fine-mapped SNP, the causal probability from FINEMAP, effect allele, effect allele frequency, sample size, corresponding gene, variant consequence, associated eQTL gene, associated eQTL associated eQTL gene and *p* value in GTEx (from <https://gtexportal.org/home/>)<sup>71</sup>. See Supplementary Table S1 for detailed information about all 25 variants and Supplementary Table S2 for fine-mapping results. \*SNP with locus-wide lowest *p* value. \*\*Odds ratio (exp(b) for logistic mixed effects model) in SHCS. \*\*\*Gene and tissue of eQTL association (*p* value of association), from GTEx: UNC5D in Esophagus\_Muscularis (*p* = 1.15243e-13), UNC5D in Esophagus\_Gastroesophageal\_Junction (*p* = 2.61983e-05).

**Genome-to-genome association analysis.** We tested for associations between each EBV variant and human SNPs. We studied 575 EBV amino acids and 52 EBV genes, for a total of 627 GWASs. The effective number of GWASs performed was 458. As covariates, we included the first six EBV principal components (51.4% deviance explained), sex, age, type 1 vs 2 of EBV (Supplementary Figure S1). The sample size ranged between 120 and 268, with a median sample size of 264. Sample size variation was due to variable missingness in the EBV data. Genomic inflation factors for each of the 627 GWASs ranged between 0.92 and 1.12.

Significant associations ( $p < 1.09 \times 10^{-10}$ ) were identified between a total of 25 human SNPs and viral variants mapping to three EBV regions (Table 2): the EBV genes **BALF5** (Fig. 2A) and **BBRF1** (Fig. 2B) and the EBV amino acid **BRLF1:p.Lys316Glu** (Fig. 2C). The minor allele frequency of all significant host SNPs was between 0.05 and 0.10. The genomic inflation factors of the three GWASs ranged between 0.95 and 0.96 (Q-Q plots shown in Supplementary Figure S5).

Strong associations were observed between 17 SNPs in the *UNC5D* region on chromosome 8 and the occurrence of very rare functional variants in the EBV **BALF5** gene (Figs. 2A, 3A), which is involved in viral DNA replication during the late phase of lytic infection. *UNC5D* is a poorly characterized gene expressed mainly in neuronal tissues, which encodes a protein that has been shown to regulate p53-dependent apoptosis in neuroblastoma cells<sup>75</sup>. The top associated SNP, rs2950922 (OR 1.31, 95% CI = 1.21–1.41,  $p = 4.2 \times 10^{-12}$ , effect allele G), is an eQTL for *UNC5D* in esophageal tissue (GTEx<sup>71</sup>).

Rare amino acid variation in **BBRF1** was found to be associated with a single SNP, rs62124869 (OR 1.29, 95% CI 1.19–1.39,  $p = 4.2 \times 10^{-11}$ , effect allele C), which maps to the non-coding RNA gene *LINC01830* (Long Intergenic Non-Protein Coding RNA 1830) on chromosome 2 (Figs. 2B, 3B).

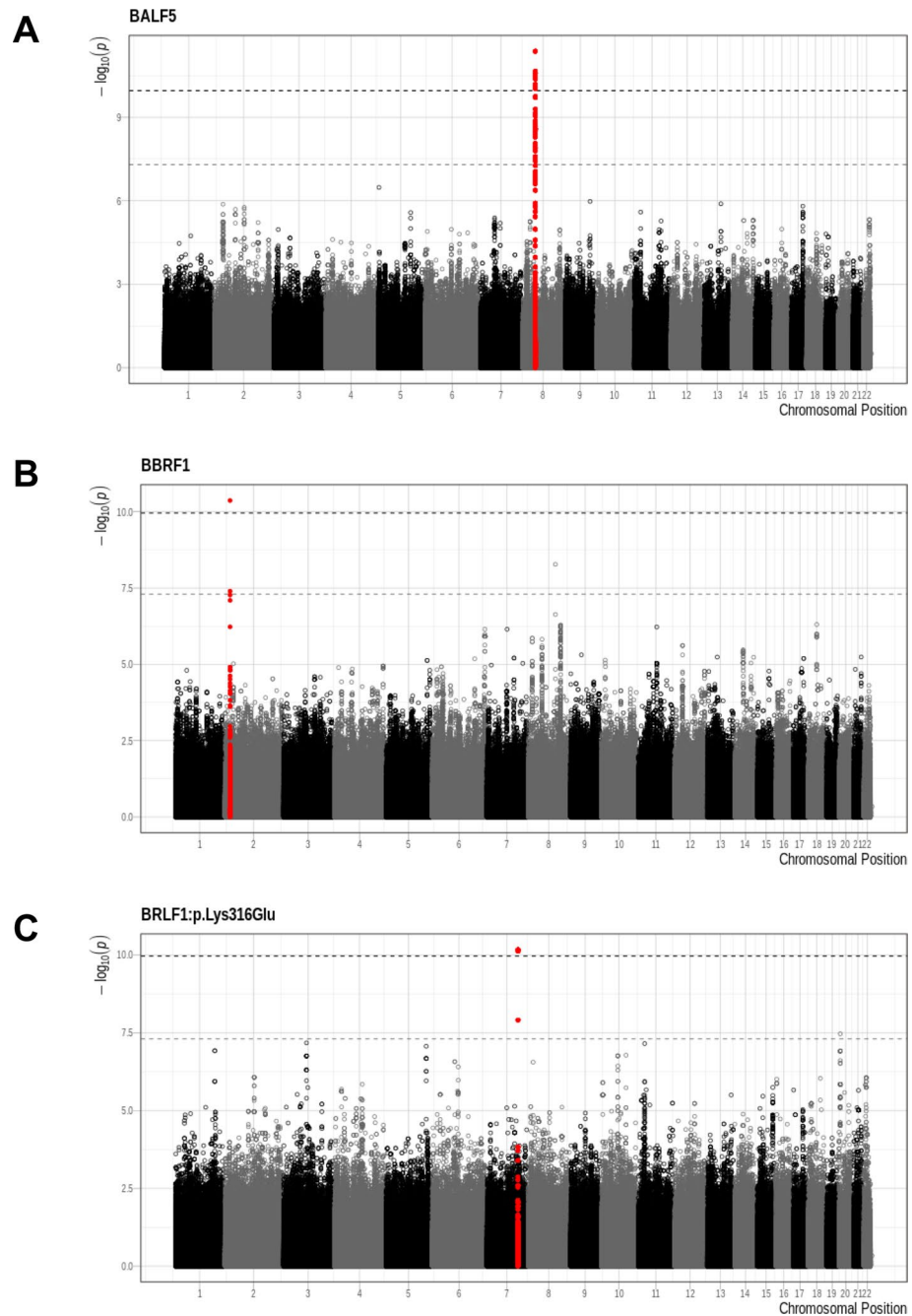
Finally, 7 SNPs mapping to a non-coding region of chromosome 7 were found to be associated with the EBV amino acid variant **BRLF1:p.Lys316Glu** (Fig. 2C, 3C). The top SNP, rs6466720, had a *p*-value of  $6.85 \times 10^{-11}$  and an OR of 1.41 (95% CI 1.28–1.58, effect allele G). **BRLF1** controls lytic reactivation of EBV from latency and regulates viral transcription. **BRLF1:p.Lys316Glu** has not been described previously, but variation at the nearby residue 377 (**BRLF1:p.Glu377Ala**) has been shown to be prevalent in cases of nasopharyngeal and gastric carcinomas in Chinese samples<sup>76</sup>. **BRLF1:p.Lys316Glu** and **BRLF1:p.Glu377Ala** are in moderate LD ( $r^2 = 0.55$ ) in our dataset.

## Discussion

Because immunosuppression—and in particular T cell deficiency—favors EBV reactivation from its latent B cell reservoir, EBV viremia is frequently detected in (untreated) HIV-infected individuals with advanced disease and low CD4+ T cell counts. We hypothesized that the reactivated virus circulating in these patients could carry sequence variants acquired during primary EBV infection, thereby providing a snapshot of early adaptation to the pressure exerted on EBV by the individual immune response.

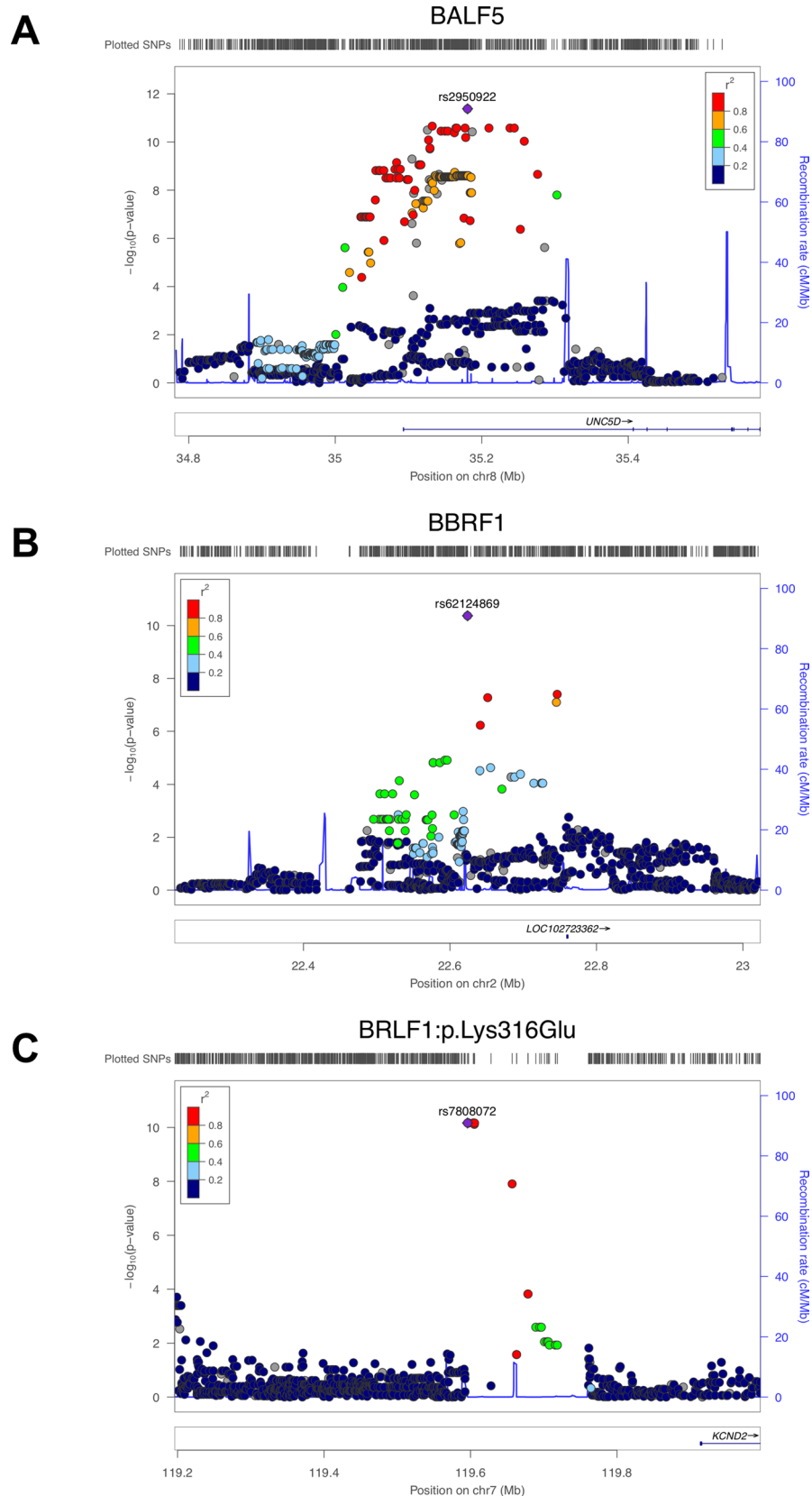
To search for a potential imprint of host genomic variation on the viral sequence, we jointly analyzed genomic information obtained from paired EBV and human samples. Viral sequence variation can be seen as an intermediate phenotype, closer to potentially causal host polymorphisms than clinically observable outcomes like viral load or disease phenotypes. As such, it allows the detection of more subtle associations, less likely to be obscured by environmental influences. In our G2G analysis, we used variation at EBV amino acid residues as outcome in multiple parallel GWAS, which allowed us to obtain effect estimations between each human genetic variant and EBV variation.

We identified two EBV genes (**BALF5**, **BBRF1**) and one EBV amino acid (BRLF1 p.Lys316Glu) as associated with three regions of the human genome, spanning altogether 25 SNPs. For the GWAS with **BALF5** as the



**Figure 2.** Significant associations—(A): *BALF5*, (B): *BBRF1*, (C): *BRLF1*:p.Lys316Glu. The x-axis represents the chromosomal position and the y-axis displays the  $-\log_{10}(p)$ . Colour alternates between chromosomes. Regions that contain statistically significant SNP are presented in red (top SNP  $\pm$  400 Kbp). The light grey dashed line represents the GWAS significance threshold of  $5 \times 10^{-8}$ , the dark grey dashed line the G2G threshold of  $1.09 \times 10^{-10}$ . This figure was produced using R<sup>65</sup>.

outcome, the associated human genomic region contains eQTLs for the nearby gene *UNC5D*, a gene shown to play a role in the regulation of apoptosis<sup>75</sup>. *BALF5* encodes the DNA polymerase catalytic subunit. Localized within replication compartments, discrete sites in nuclei, it is one of the six proteins forming the viral replication complex along with processivity factor, primase, primase-associated factor, helicase, and ssDNA-binding protein. It contributes to the replication of viral genomic DNA in the late phase of lytic infection, producing long concatemeric DNA<sup>77</sup>. *BRLF1* encodes an immediate-early transcriptional activator. It induces the initiation of viral lytic gene expression and lytic reactivation from latency, a key process in parrying the host immune response. It has also been shown to both upregulate human *TNFRSF6B* by directly binding to its receptor and to interact with human transcription factor protein *ATF7IP*, potentially regulating host genes in virus-infected cells<sup>78</sup>. EBV capsid assembly includes the portal oligomer, encoded by *BBRF1*, through which viral DNA is



**Figure 3.** Locuszoom plots. Locuszoom plots for the three EBV association signals highlighted in red in Fig. 2 (A: *BALF5*, B: *BBRF1*, C: *BRLF1*:p.Lys316Glu). This Figure was produced using LocusZoom<sup>74</sup>.



translocated during DNA packaging. Forming a homododecamer, it is translocated to the nucleus by viral scaffold protein and binds to terminase, a molecular motor, that translocates the viral DNA. The deletion of **BBRF1** may be used in the production of DNA-free virus-like particles/light particles for preventive vaccines against hepatitis B and human papillomaviruses<sup>79</sup>.

Our study is limited by its small sample size and by the complexity of correcting for human and EBV population stratification. Indeed, if not carefully controlled for, the existence of population structure in the host and pathogen genome might create spurious associations or decrease real signals in G2G analyses, resulting in both type I and type II errors. With a mixed model approach and the inclusion of pathogen principal components as covariates, the genomic inflation factors of our GWAS ranged between 0.92 and 1.12. This wide range of genomic inflation factors is likely due to a combination of small sample size and complex statistical model. To prevent false positives, we adjusted for genomic inflation when extracting significant SNPs and used a conservative G2G significance threshold of  $5 \times 10^{-8}$  divided by the effective number of GWAS performed. Although viral genetic variation is a more precise phenotype to study than traditional outcomes, it comes at the price of decreased power due to the high-dimensional outcome. The significance threshold is thus much lower than in a single GWAS. To limit the number of statistical tests performed, we restricted our analysis to common gene and amino acid variation.

Our analyses have been performed using historical samples collected from untreated HIV-infected individuals. Considering the natural history of EBV infection in humans and its high likelihood to be acquired during the first 2 decades of life, we postulate that intra-host adaptation of EBV happened before HIV infection, i.e. with a normally functioning immune system. At the time of sample collection, all study participants had advanced immunosuppression with low CD4 + T cell counts ( $< 200$  cells/mm<sup>3</sup> of blood). We therefore assume an absence of selective pressure on EBV at that time. These assumptions limit obviously the generalizability of our findings to non-HIV-infected population. Similar studies performed during primary EBV infection or in other specific population (e.g. bone-marrow transplant recipients) would help better delineate the global impact of intra-host selection on EBV sequence variation.

Our study provides a preliminary list of statistical associations between the EBV and the human genomes. The cataloguing of the sites of host–pathogen genomic conflict is potentially useful for further functional exploration, as has been demonstrated for HIV and HCV infections. Our results require replication and validation in different cohorts and settings. Importantly, larger sample sizes will be needed to increase power and provide more robust estimations.

### Data availability

The datasets generated during and/or analysed during the current study are available in the following Zenodo repositories: G2G results are in <https://doi.org/10.5281/zenodo.4289138>, pathogen data in <https://doi.org/10.5281/zenodo.4011995>.

### Code availability

EBV data preparation: [https://gitlab.com/e2lab/vir\\_var\\_calling](https://gitlab.com/e2lab/vir_var_calling). G2G Analysis: <https://github.com/sinarueeger/G2G-EBV-manuscript>.

Received: 2 December 2020; Accepted: 11 February 2021

Published online: 25 February 2021

### References

- Chapman, S. J. & Hill, A. V. S. Human genetic susceptibility to infectious disease. *Nat. Rev. Genet.* **13**, 175–188 (2012).
- Casanova, J.-L. & Abel, L. The human genetic determinism of life-threatening infectious diseases: genetic heterogeneity and physiological homogeneity?. *Hum. Genet.* **139**, 681–694 (2020).
- Visscher, P. M. *et al.* 10 years of GWAS discovery: biology, function, and translation. *Am. J. Hum. Genet.* **101**, 5–22 (2017).
- Timmann, C. *et al.* Genome-wide association study indicates two novel resistance loci for severe malaria. *Nature* **489**, 443–446 (2012).
- McLaren, P. J. *et al.* Association study of common genetic variants and HIV-1 acquisition in 6,300 infected cases and 7,200 controls. *PLoS Pathog.* **9**, e1003515 (2013).
- McLaren, P. J. *et al.* Polymorphisms of large effect explain the majority of the host genetic contribution to variation of HIV-1 virus load. *PNAS* **112**, 14658–14663 (2015).
- Rubicz, R. *et al.* A genome-wide integrative genomic study localizes genetic factors influencing antibodies against Epstein–Barr virus nuclear antigen 1 (EBNA-1). *PLoS Genet.* **9**, e1003147 (2013).
- Zhou, Y. *et al.* Genetic loci for Epstein–Barr virus nuclear antigen-1 are associated with risk of multiple sclerosis. *Mult. Scler.* **22**, 1655–1664 (2016).
- Hammer, C. *et al.* Amino acid variation in HLA class II proteins is a major determinant of humoral response to common viruses. *Am. J. Hum. Genet.* **97**, 738–743 (2015).
- Ge, D. *et al.* Genetic variation in *IL28B* predicts hepatitis C treatment-induced viral clearance. *Nature* **461**, 399–401 (2009).
- Alizon, S., Luciani, F. & Regoes, R. R. Epidemiological and clinical consequences of within-host evolution. *Trends Microbiol.* **19**, 24–32 (2011).
- Fraser, C. *et al.* Virulence and pathogenesis of HIV-1 infection: an evolutionary perspective. *Science* **343**, 1243727 (2014).
- Farci, P. *et al.* The outcome of acute hepatitis C predicted by the evolution of the viral quasispecies. *Science* **288**, 339–344 (2000).
- Didelot, X., Walker, A. S., Peto, T. E., Crook, D. W. & Wilson, D. J. Within-host evolution of bacterial pathogens. *Nat. Rev. Microbiol.* **14**, 150–162 (2016).
- Bartha, I. *et al.* A genome-to-genome analysis of associations between human genetic variation, HIV-1 sequence diversity, and viral control. *Elife* **2**, e01123 (2013).
- Cohen, J. I. Epstein–Barr virus vaccines. *Clin. Transl. Immunol.* **4**, e32 (2015).
- Ansari, M. A. *et al.* Genome-to-genome analysis highlights the effect of the human innate and adaptive immune systems on the hepatitis C virus. *Nat. Genet.* **49**, 666–673 (2017).
- Ansari, M. A. *et al.* Interferon lambda 4 impacts the genetic diversity of hepatitis C virus. *Elife* **8**, e42463 (2019).

19. Chaturvedi, N. *et al.* Adaptation of hepatitis C virus to interferon lambda polymorphism across multiple viral genotypes. *eLife* **8**, e42542 (2019).
20. Duffy, S., Shackelton, L. A. & Holmes, E. C. Rates of evolutionary change in viruses: patterns and determinants. *Nat. Rev. Genet.* **9**, 267–276 (2008).
21. Cudini, J. *et al.* Human cytomegalovirus haplotype reconstruction reveals high diversity due to superinfection and evidence of within-host recombination. *Proc. Natl. Acad. Sci. USA* **116**, 5693 (2019).
22. Kwok, H. *et al.* Genomic diversity of Epstein–Barr virus genomes isolated from primary nasopharyngeal carcinoma biopsy samples. *J. Virol.* **88**, 10662–10672 (2014).
23. Palsler, A. L. *et al.* Genome diversity of Epstein–Barr virus from multiple tumor types and normal infection. *J. Virol.* **89**, 5222–5237 (2015).
24. Balfour, H. H. *et al.* Age-specific prevalence of Epstein–Barr virus infection among individuals aged 6–19 years in the United States and factors affecting its acquisition. *J. Infect. Dis.* **208**, 1286–1293 (2013).
25. Green, M. & Michaels, M. G. Epstein–Barr virus infection and posttransplant lymphoproliferative disorder: EBV and PTLID. *Am. J. Transplant.* **13**, 41–54 (2013).
26. Pender, M. P. The essential role of Epstein–Barr virus in the pathogenesis of multiple sclerosis. *Neuroscientist* **17**, 351–367 (2011).
27. Pender, M. P. & Burrows, S. R. Epstein–Barr virus and multiple sclerosis: potential opportunities for immunotherapy. *Clin. Transl. Immunol.* **3**, e27 (2014).
28. Farina, A. *et al.* Epstein–Barr virus lytic infection promotes activation of Toll-like receptor 8 innate immune response in systemic sclerosis monocytes. *Arthritis Res. Ther.* **19**, 39 (2017).
29. Ruprecht, K. The role of Epstein–Barr virus in the etiology of multiple sclerosis: a current review. *Expert Rev. Clin. Immunol.* **0**, 1–15 (2020).
30. Young, L. S. & Rickinson, A. B. Epstein–Barr virus: 40 years on. *Nat. Rev. Cancer* **4**, 757–768 (2004).
31. Ko, Y.-H. EBV and human cancer. *Exp. Mol. Med.* **47**, e130 (2015).
32. de Martel, C. *et al.* Global burden of cancers attributable to infections in 2008: a review and synthetic analysis. *Lancet Oncol.* **13**, 607–615 (2012).
33. Khan, G. & Hashim, M. J. Global burden of deaths from Epstein–Barr virus attributable malignancies 1990–2010. *Infect. Agents Cancer* **9**, 38 (2014).
34. Hammerschmidt, W. & Sugden, B. Replication of Epstein–Barr viral DNA. *Cold Spring Harb. Perspect. Biol.* **5**, a013029 (2013).
35. Xu, M. *et al.* Genome sequencing analysis identifies Epstein–Barr virus subtypes associated with high risk of nasopharyngeal carcinoma. *Nat. Genet.* <https://doi.org/10.1038/s41588-019-0436-5> (2019).
36. Wegner, F., Lassalle, F., Depledge, D. P., Balloux, F. & Breuer, J. Co-evolution of sites under immune selection shapes Epstein–Barr Virus population structure. *Mol. Biol. Evol.* <https://doi.org/10.1093/molbev/msz152> (2019).
37. Kaymaz, Y. *et al.* Epstein Barr virus genomes reveal population structure and type 1 association with endemic Burkitt lymphoma. *bioRxiv* <https://doi.org/10.1101/689216> (2019).
38. Chiara, M. *et al.* Geographic population structure in Epstein–Barr virus revealed by comparative genomics. *Genome Biol. Evol.* **8**, 3284–3291 (2016).
39. Hayward, T. A. *et al.* Antibody response to common human viruses is shaped by genetic factors. *J. Allergy Clin. Immunol.* **143**, 1640–1643 (2019).
40. The Swiss HIV Cohort Study *et al.* Cohort profile: the Swiss HIV Cohort Study. *Int. J. Epidemiol.* **39**, 1179–1189 (2010).
41. Depledge, D. P. *et al.* Specific capture and whole-genome sequencing of viruses from clinical samples. *PLoS ONE* **6**, e27805 (2011).
42. Bolger, A. M., Lohse, M. & Usadel, B. Trimmomatic: a flexible trimmer for Illumina sequence data. *Bioinformatics* **30**, 2114–2120 (2014).
43. Lassmann, T. TagDust2: a generic method to extract reads from sequencing data. *BMC Bioinform.* **16**, 24 (2015).
44. Fu, L., Niu, B., Zhu, Z., Wu, S. & Li, W. CD-HIT: accelerated for clustering the next-generation sequencing data. *Bioinformatics* **28**, 3150–3152 (2012).
45. DePristo, M. A. *et al.* A framework for variation discovery and genotyping using next-generation DNA sequencing data. *Nat. Genet.* **43**, 491–498 (2011).
46. Van der Auwera, G. A. *et al.* From FastQ data to high confidence variant calls: the Genome Analysis Toolkit best practices pipeline. *Curr. Protoc. Bioinform.* **43**, 11.10.1–11.10.33 (2013).
47. Li, H. Aligning sequence reads, clone sequences and assembly contigs with BWA-MEM. arXiv:1303.3997 [q-bio] (2013).
48. McKenna, A. *et al.* The genome analysis toolkit: a MapReduce framework for analyzing next-generation DNA sequencing data. *Genome Res.* **20**, 1297–1303 (2010).
49. Correia, S. *et al.* Natural variation of Epstein–Barr virus genes, proteins, and primary microRNA. *J. Virol.* **91**(15), e00375–17 (2017).
50. Wei, Z., Wang, W., Hu, P., Lyon, G. J. & Hakonarson, H. SNVer: a statistical tool for variant calling in analysis of pooled or individual next-generation sequencing data. *Nucleic Acids Res.* **39**, e132 (2011).
51. Koboldt, D. C. *et al.* VarScan: variant detection in massively parallel sequencing of individual and pooled samples. *Bioinformatics* **25**, 2283–2285 (2009).
52. Cingolani, P. *et al.* A program for annotating and predicting the effects of single nucleotide polymorphisms, SnpEff: SNPs in the genome of *Drosophila melanogaster* strain w1118; iso-2; iso-3. *Fly (Austin)* **6**, 80–92 (2012).
53. Köster, J. & Rahmann, S. Snakemake—a scalable bioinformatics workflow engine. *Bioinformatics* **28**, 2520–2522 (2012).
54. *Python Language Reference*. (Python Software Foundation).
55. Brinda, K., Boeva, V. & Kucherov, G. RNF: a general framework to evaluate NGS read mappers. *Bioinformatics* **32**, 136–139 (2015).
56. Josse, J. & Husson, F. missMDA: A package for handling missing values in multivariate data analysis. *J. Stat. Softw.* **70**, 1–31 (2016).
57. Loh, P.-R. *et al.* Reference-based phasing using the Haplotype Reference Consortium panel. *Nat. Genet.* **48**, 1443–1448 (2016).
58. Durbin, R. Efficient haplotype matching and storage using the positional Burrows–Wheeler transform (PBWT). *Bioinformatics* **30**, 1266–1272 (2014).
59. The 1000 Genomes Project Consortium. A global reference for human genetic variation. *Nature* **526**, 68–74 (2015).
60. Chang, C. C. *et al.* Second-generation PLINK: rising to the challenge of larger and richer datasets. *GigaScience* **4**, 7 (2015).
61. Yang, J., Lee, S. H., Goddard, M. E. & Visscher, P. M. GCTA: a tool for genome-wide complex trait analysis. *Am. J. Hum. Genet.* **88**, 76–82 (2011).
62. Yang, J., Zaitlen, N. A., Goddard, M. E., Visscher, P. M. & Price, A. L. Advantages and pitfalls in the application of mixed-model association methods. *Nat. Genet.* **46**, 100–106 (2014).
63. Naret, O. *et al.* Correcting for population stratification reduces false positive and false negative results in joint analyses of host and pathogen genomes. *Front. Genet.* **9**, 266 (2018).
64. Landgraf, A. J. & Lee, Y. Dimensionality reduction for binary data through the projection of natural parameters. arXiv:1510.06112 [stat] (2015).
65. R Core Team. *R: A Language and Environment for Statistical Computing*. (R Foundation for Statistical Computing, 2020).
66. Gao, X., Starmer, J. & Martin, E. R. A multiple testing correction method for genetic association studies using correlated single nucleotide polymorphisms. *Genet. Epidemiol.* **32**, 361–369 (2008).
67. Benner, C. *et al.* FINEMAP: efficient variable selection using summary data from genome-wide association studies. *Bioinformatics* **32**, 1493–1501 (2016).

68. Benner, C. *et al.* Prospects of fine-mapping trait-associated genomic regions by using summary statistics from genome-wide association studies. *Am. J. Hum. Genet.* **101**, 539–551 (2017).
69. Vösa, U. *et al.* Unraveling the polygenic architecture of complex traits using blood eQTL meta-analysis. <https://doi.org/10.1101/447367> (2018).
70. Ferreira, M. A. R. *et al.* Gene-based analysis of regulatory variants identifies 4 putative novel asthma risk genes related to nucleotide synthesis and signaling. *J. Allergy Clin. Immunol.* **139**, 1148–1157 (2017).
71. Carithers, L. J. *et al.* A novel approach to high-quality postmortem tissue procurement: the GTEx project. *Biopreserv. Biobank* **13**, 311–319 (2015).
72. Li, H. A statistical framework for SNP calling, mutation discovery, association mapping and population genetical parameter estimation from sequencing data. *Bioinformatics* **27**, 2987–2993 (2011).
73. Sandmann, S. *et al.* Evaluating variant calling tools for non-matched next-generation sequencing data. *Sci. Rep.* **7**, 43169 (2017).
74. Pruim, R. J. *et al.* LocusZoom: regional visualization of genome-wide association scan results. *Bioinformatics* **26**, 2336–2337 (2010).
75. Wang, H. *et al.* Unc5D regulates p53-dependent apoptosis in neuroblastoma cells. *Mol. Med. Rep.* **9**, 2411–2416 (2014).
76. Jia, Y. *et al.* Sequence analysis of the Epstein–Barr virus (EBV) BRLF1 gene in nasopharyngeal and gastric carcinomas. *Viol. J.* **7**, 341 (2010).
77. Kiehl, A. & Dorsky, D. I. Cooperation of EBV DNA polymerase and EA-D(BMRF1) in vitro and colocalization in nuclei of infected cells. *Virology* **184**, 330–340 (1991).
78. Darr, C. D., Mauser, A. & Kenney, S. Epstein–Barr virus immediate-early protein BRLF1 induces the lytic form of viral replication through a mechanism involving phosphatidylinositol-3 kinase activation. *J. Virol.* **75**, 6135–6142 (2001).
79. Pavlova, S. *et al.* An Epstein–Barr virus mutant produces immunogenic defective particles devoid of viral DNA. *J. Virol.* **87**, 2011–2022 (2013).

## Acknowledgements

This study was supported by the Leenaards Foundation (Leenaards Prize 2015 to JF and EZ). This study has also been partly financed within the framework of the Swiss HIV Cohort Study, supported by the Swiss National Science Foundation (Grant #177499), by SHCS Project #743 and by the SHCS research foundation. The data are gathered by the Five Swiss University Hospitals, two Cantonal Hospitals, 15 affiliated hospitals and 36 private physicians (listed in <http://www.shcs.ch/180-health-care-providers>).

## Author contributions

E.Z. and J.F. conceived and supervised the work. S.R. and C.H. performed the association analyses. D.P.D., S.M. and J.B. performed EBV sequencing and curated the data. A.L. prepared and analysed the viral sequencing data. P.J.M., D.L. and O.N. contributed to the design of the work. N.K., E.B., M.C., H.F.G., C.R.K. and A.R. recruited the study participants and collected the samples and associated data. S.R., C.H., A.L. and J.F. wrote the paper. All authors reviewed the manuscript and approved the submission.

## Competing interests

CH is an employee of Genentech.

## Additional information

**Supplementary Information** The online version contains supplementary material available at <https://doi.org/10.1038/s41598-021-84070-7>.

**Correspondence** and requests for materials should be addressed to J.F.

**Reprints and permissions information** is available at [www.nature.com/reprints](http://www.nature.com/reprints).

**Publisher's note** Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.



**Open Access** This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons licence, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons licence and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this licence, visit <http://creativecommons.org/licenses/by/4.0/>.

© The Author(s) 2021

## the Swiss HIV Cohort Study

Karoline Aebi-Popp<sup>13</sup>, Alexia Anagnostopoulos<sup>10</sup>, Manuel Battegay<sup>16</sup>, Enos Bernasconi<sup>8</sup>, Jürg Böni<sup>17</sup>, Dominique Braun<sup>10</sup>, Heiner Bucher<sup>18</sup>, Alexandra Calmy<sup>19</sup>, Matthias Cavassini<sup>9</sup>, Angela Ciuffi<sup>20</sup>, Guenter Dollenmaier<sup>21</sup>, Matthias Egger<sup>22</sup>, Luigia Elzi<sup>16</sup>, Jan Fehr<sup>10</sup>, Jacques Fellay<sup>1,2,15</sup>, Hansjakob Furrer<sup>13</sup>, Christoph Fux<sup>23</sup>, Huldrych F. Günthard<sup>10</sup>, David Haerry<sup>24</sup>, Barbara Hasse<sup>10</sup>, Hans Hirsch<sup>25,26</sup>, Matthias Hoffmann<sup>11</sup>, Irene Hösli<sup>27</sup>, Michael Huber<sup>17</sup>, Christian R. Kahlert<sup>11,12</sup>, Laurent Kaiser<sup>28</sup>, Olivia Keiser<sup>29</sup>, Thomas Klimkait<sup>25</sup>, Lisa Kottanattu<sup>30</sup>, Roger Kouyos<sup>10</sup>, Helen Kovari<sup>10</sup>, Bruno Ledergerber<sup>10</sup>,

**Gladys Martinetti<sup>31</sup>, Begoña Martinez de Tejada<sup>32</sup>, Catia Marzolini<sup>16</sup>, Karin Metzner<sup>10</sup>, Nicolas Müller<sup>10</sup>, Dunja Nicca<sup>11</sup>, Paolo Paioni<sup>33</sup>, Giuseppe Pantaleo<sup>34</sup>, Matthieu Perreau<sup>34</sup>, Andri Rauch<sup>13</sup>, Christoph Rudin<sup>35</sup>, Alexandra Scherrer<sup>10,36</sup>, Patrick Schmid<sup>11</sup>, Roberto Speck<sup>10</sup>, Marcel Stöckle<sup>16</sup>, Philip Tarr<sup>37</sup>, Alexandra Trkola<sup>17</sup>, Pietro Vernazza<sup>11</sup>, Noémie Wagner<sup>38</sup>, Gilles Wandeler<sup>13</sup>, Rainer Weber<sup>10</sup> & Sabine Yerly<sup>39</sup>**

<sup>16</sup> Division of Infectious Diseases and Hospital Epidemiology, University Hospital Basel, University of Basel, Basel, Switzerland. <sup>17</sup>Institute of Medical Virology, University of Zürich, Zurich, Switzerland. <sup>18</sup>Basel Institute for Clinical Epidemiology and Biostatistics, University Hospital Basel, University of Basel, Basel, Switzerland. <sup>19</sup>Division of Infectious Diseases, University Hospital Geneva, University of Geneva, Geneva, Switzerland. <sup>20</sup>Institute of Microbiology, University Hospital Lausanne, University of Lausanne, Lausanne, Switzerland. <sup>21</sup>Centre for Laboratory Medicine, St. Gallen, Canton St. Gallen, Switzerland. <sup>22</sup>Institute of Social and Preventive Medicine, University of Bern, Bern, Switzerland. <sup>23</sup>Clinic for Infectious Diseases and Hospital Hygiene, Kantonsspital Aarau, Aarau, Switzerland. <sup>24</sup>Deputy of the Patient Organization "Positive Council", Zurich, Switzerland. <sup>25</sup>Division Infection Diagnostics, Department Biomedicine - Petersplatz, University of Basel, Basel, Switzerland. <sup>26</sup>Division of Infectious Diseases and Hospital Epidemiology, University Hospital Basel, Basel, Switzerland. <sup>27</sup>Clinic for Obstetrics, University Hospital Basel, University of Basel, Basel, Switzerland. <sup>28</sup>Division of Infectious Diseases and Laboratory of Virology, University Hospital Geneva, University of Geneva, Geneva, Switzerland. <sup>29</sup>Institute of Global Health, University of Geneva, Geneva, Switzerland. <sup>30</sup>Ospedale Regionale di Bellinzona e Valli, Bellinzona, Switzerland. <sup>31</sup>Cantonal Institute of Microbiology, Bellinzona, Switzerland. <sup>32</sup>Department of Obstetrics and Gynecology, University Hospital Geneva, University of Geneva, Geneva, Switzerland. <sup>33</sup>University Children's Hospital, University of Zurich, Zurich, Switzerland. <sup>34</sup>Division of Immunology and Allergy, University Hospital Lausanne, University of Lausanne, Lausanne, Switzerland. <sup>35</sup>University Childrens Hospital, University of Basel, Basel, Switzerland. <sup>36</sup>Swiss HIV Cohort Study, Data Centre, Zurich, Switzerland. <sup>37</sup>Kantonsspital Baselland, University of Basel, Basel, Switzerland. <sup>38</sup>Pediatrics, University Hospital Geneva, University of Geneva, Geneva, Switzerland. <sup>39</sup>Laboratory of Virology, University Hospital Geneva, University of Geneva, Geneva, Switzerland.