

15

Abstract

16 Genome-wide association studies (GWAS) have considerably advanced our
17 understanding of human traits and diseases. With the increasing availability of whole genome
18 sequences (WGS) for pathogens, it is important to establish whether GWAS of viral genomes
19 could reveal important biological insights. Here we perform the first proof of concept viral
20 GWAS examining drug resistance (DR), a phenotype with well understood genetics.

21 We performed a GWAS of DR in a sample of 343 HIV subtype C patients failing 1st
22 line antiretroviral treatment in rural KwaZulu-Natal, South Africa. The majority and minority
23 variants within each sequence were called using PILON, and GWAS was performed within
24 PLINK. HIV WGS from patients failing on different antiretroviral treatments were compared
25 to sequences derived from individuals naive to the respective treatment.

26 GWAS methodology was validated by identifying five associations on a genetic level
27 that led to amino acid changes known to cause DR. Further, we highlighted the ability of
28 GWAS to identify epistatic effects, identifying two replicable variants within amino acid 68
29 of the reverse transcriptase protein previously described as potential fitness compensatory
30 mutations. A possible additional DR variant within amino acid 91 of the matrix region of the
31 Gag protein was associated with tenofovir failure, highlighting the ability of GWAS to
32 identify variants outside classical candidate genes. Our results also suggest a polygenic
33 component to DR.

34 These results validate the applicability of GWAS to HIV WGS data even in relative
35 small samples, and emphasise how high throughput sequencing can provide novel and
36 clinically relevant insights. Further they suggested that for viruses like HIV, population
37 structure was only minor concern compared to that seen in bacteria or parasite GWAS. Given

38 the small genome length and reduced burden for multiple testing, this makes HIV an ideal
39 candidate for GWAS.

40

41

Introduction

42 Genome-wide association studies (GWAS) have led to significant advances in the
43 understanding of complex human traits and diseases. They involve the analysis of hundreds
44 of thousands or millions of common genetic variants, usually single nucleotide
45 polymorphisms (SNPs), testing for an association between each variant and a phenotype (see
46 [1]). This allows for the analysis of many variants across the genome, blind to their location
47 or functionality. This approach has identified hundreds of causal risk variants for dozens of
48 diseases in the last decade (e.g. [2-4]), each a potential drug target for novel treatments.
49 These advances were made possible due to the availability of cost effective SNP genotyping
50 technology which capture known common genetic variants. The limitation of this approach is
51 that it misses variants absent from the chip, especially rare or *de novo* mutations. For this
52 reason, genetic research is increasingly moving towards whole genome sequencing
53 approaches to capture the full range of genetic variants in a population.

54 In this respect, the field of pathogen genomics is quickly catching up with human
55 genomics, with international collaborations currently generating thousands of whole genome
56 sequences (WGS) for pathogens such as HIV and malaria (e.g. the PANGEA Consortium[5]
57 and the MalariaGen Consortium[6]). These WGS allow for the application of GWAS-style
58 identification of novel genetic risk variants without the need for SNP genotyping chips.

59 A GWAS approach has previously been successfully applied to other non-virus
60 pathogen, almost always using treatment resistance or failure as the phenotype[7]. These

61 studies have included *Plasmodium falciparum*[8], *Mycobacterium tuberculosis*[9],
62 *Staphylococcus aureus*[10] and *Streptococcus pneumoniae*[11]. Sample sizes have ranged
63 from 75 to 3,701 sequences, and in even smaller samples have identified both novel and
64 known variants that capture almost all the variation in treatment outcome.

65

66 However it is still unclear how well suited the viral genome is to a GWAS approach.
67 The only viral GWAS to date combined GWAS of human SNP and HIV amino acid data, and
68 identified multiple host genetic variants in the HLA region associated with HIV amino acid
69 diversity[12]. However they found no associations between the HIV genome and their
70 outcome of interest, viral load. The high percentage of coding sequence in viral genomes and
71 overlapping reading frames may constrain the polygenic architecture for which GWAS was
72 conceived: with many variants each of individually small effect. Another limitation of
73 previous studies was that they did not allow for heterozygosity. Heterozygosity at a locus can
74 arise due to mixed infections or within-host pathogen genetic diversification. Although this is
75 rare in most pathogens studied with GWAS to date, it is highly relevant to many viral
76 infections. Lastly, parasite and bacterial GWAS have observed a large level of population
77 structure presumably due in part to homologous recombination and recent selection[13].
78 Given the challenges faced by previous analyses, more work is needed to properly define the
79 genomic architecture of viruses and whether it is suitable to a GWAS style approach.

80 To validate the effectiveness of a viral GWAS we aimed to replicate the success of
81 bacterial GWAS and focus on a phenotype known to be under strong selection pressure,
82 specifically antiretroviral therapy (ART) resistance in HIV. The provision of ART to over 6.2
83 million people living with HIV in sub-Saharan Africa has been one of the most successful
84 public health interventions ever undertaken[14], improving life expectancy[15], and reducing
85 transmission[16, 17]. As a result, ART has been one of the most potent selection pressures on

86 HIV. Given its importance to global health, resistance to ART has been well studied in HIV
87 with many amino acid changes known to lead to DR [18]. Thus, DR is a useful phenotype
88 for validating GWAS in HIV as findings can be compared to the existing literature as well as
89 to large publically available databases of genes involved in HIV DR. In this study, we aim to
90 identify known variants and validate the applicability of GWAS methods to the HIV genome.

91 **Results**

92 ***Genomic architecture of HIV SNPs***

93 343 samples with phenotype and genotype data remained after variant calling and quality
94 control (Table 1). A total of the 5379 SNPs with a minor allele frequency $\geq 1\%$ were
95 identified. An excess of rare variants was observed with a mean allele frequency of 11.3%
96 and median of 6.0% (see Supplementary Figure 1). Additionally 2502 variants were
97 identified with a frequency less than 1% though not included in the analyses. Variants were
98 evenly distributed across the genome, despite missingness differing by region (see
99 Supplementary Figure 2). The permuted threshold for genome-wide significance was $p=7E-5$,
100 less stringent than that derived by Bonferroni correction for the number of variants ($p=9.3E-$
101 6) and suggesting that there was substantial correlation between SNPs. This correlation is
102 expected, due to the close proximity of SNPs in WGS data which leads to linkage
103 disequilibrium and the non-independence of tests. As such, genome-wide significance was
104 determined using the permutation adjusted p- permutation adjusted p-value threshold. SNPs
105 were labelled by their base pair position plus reference allele, e.g. 1A. SNPs were also linked
106 to their corresponding amino acid position in the different HIV proteins using reference
107 sequence AF411967.

108 ***Validating GWAS with known DR variants***

109 GWAS was performed to identify variants associated with drug resistance. The drug
110 resistance phenotype was binary for each drug and defined as any history, or not, of failure
111 while treated with the given drug. Failure was defined as at least one measure of viral load
112 >1000 copies/ml after 12 months of treatment. GWAS identified eight independent
113 associations at permutation adjusted genome-wide significance. Five of the associations were
114 known loci involved in DR and all but one were in the reverse transcriptase region (RT), the
115 functional target of these drugs (see Table 2). Failure on tenofovir was associated with three
116 known SNPs (2730G, 2852A and 2880T, see Figure 1) in the RT region, at amino acid
117 positions 65, 106, and 115, of which position 65 and 115 were known tenofovir DR variants
118 and position 106 was previously associated with DR with the most common drugs used in
119 combination with tenofovir. Treatment with zidovudine was associated with SNP 2745G, a
120 known drug variant in RT amino acid 70 (Supplementary Figure 3). Nevirapine treatment
121 was associated with a SNP (3078G) at RT 181, again a previously known DR variant
122 (Supplementary Figure 4). No associations were seen with known resistance variants for
123 lopinavir, efavirenz and stavudine (Supplementary Figure 5, 6 and 7). These results remained
124 significant after correction for confounding from population structure and length of treatment
125 (Supplementary Table 1)

126 While our analyses identified several known variants for DR, not all were identified.
127 However, it is well known in GWAS studies that sample size is a critical limitation, with
128 additional SNPs identified when larger samples are available. We observed a weak positive
129 correlation in our analyses between the number of significant associations per drug and
130 sample size ($R^2=22\%$). Looking at known DR mutations[18] with data available showed an
131 excess of significant associations compared to expectation by chance, with 12% containing a
132 variant at genome-wide significance and a further 41% containing at least one at nominal

133 significance, despite incomplete coverage ($p < 0.001$; see Supplementary Table 2). This trend
134 was especially clear within the primary resistance mutations.

135 *Identification of novel variants*

136 As well as known drug resistance variants, additional associations were observed. The
137 first was two associations within RT amino acid 68. The first was between tenofovir failure
138 and SNP 2738G resulting in a change from serine to glycine. Replication was performed
139 using the Stanford University HIV Drug Resistance Database[19, 20]. Subtype C sequences
140 within the Stanford database from individuals failing to tenofovir or other NRTIs ($n=9,357$)
141 all had the reference (serine) amino acid. For sequences showing resistance to tenofovir,
142 however, 5.5% had glycine at this position ($n=488$, $p < 0.0001$ compared to non exposed
143 distribution). Stavudine resistance was showed associations with a different SNP (2739A)
144 from serine to asparagine (Supplementary Figure 3). However further investigation showed
145 this to be an association with the negatively correlated drug tenofovir which had a p-value
146 just below genome-wide significance for this SNP ($p=7.1E-4$). This was clear both from the
147 fact the reference sequence allele was associated with stavudine DR, and from the results of
148 the replication. For sequences failing on tenofovir 4.7% had asparagine at position 68
149 ($p < 0.0001$), while for sequences failing on stavudine ($n=2,800$) no asparagine variants were
150 observed. While not known drug resistance variants, amino acid 68 (specifically the change
151 to glycine) has been suggested as a compensatory mutation for reduced fitness due to the
152 drug resistance variant in amino acid 65[21, 22]. Indeed epistasis was observed between the
153 significant SNPs in amino acid 68 and those in amino acids 65 and 106 (Supplementary
154 Table 3).

155 For tenofovir failure an association was also seen with SNP 1063A in amino acid 91
156 of the matrix region, an entirely novel association (see Figure 2). Whilst not available in the

157 Stanford database, we compared our results to the Los Alamos HIV Sequence Database drug
158 naïve WGS at the amino acid 91 of the matrix region. Interestingly, for the amino acid 91 we
159 observed a high level of genetic variation, with coding for nine amino acids. Focusing on the
160 associated genetic change, we observed significantly different ($p < 0.0001$) frequencies in the
161 drug naïve sample (37% G vs. 61% A) compared to the tenofovir-exposed sequences in our
162 sample (65% G vs. 35% A). Our WGS tenofovir naïve cases had a same frequency as the
163 publically available sequences (37% G). While not an independent replication, this lends
164 some support to our finding.

165 *Population stratification and cryptic relatedness*

166 A concern in GWAS is the possibility of confounding by population stratification,
167 which can lead to a systematic inflation in the number of false positives. QQ plots are a
168 standard tool for testing for inflation in GWAS, plotting observed p-values across the genome
169 compared to expected p-value distribution. These suggested a systematic deflation in p-values
170 in this study, with genomic lambdas between 0.66-0.80. The lambda value is derived from the
171 median observed chi squared statistic divided by the median expected chi squared statistic
172 (for $p=0.5$). Under the null distribution, a lambda of 1 is expected, with a value above 1.05-
173 1.10 usually taken as evidence of inflation. However, the excess of very rare variants (see
174 Supplementary Figure 1) prevented a normal distribution of p-values, with a reduced number
175 of significant SNPs compared to expected under the null. Restricting the analysis to SNPs
176 where minor allele frequency was at least 10% supported this hypothesis, with an increase in
177 genomic lambdas (0.81-1.00). To account for this, we compared our distribution of p-values
178 to those when the phenotype was permuted within our data. This removed the systematic
179 deflation in our expected vs. observed p-value distributions (lambdas 0.99-1.36, median
180 1.076), now showing a distribution close to null for the majority of SNPs (see Supplementary
181 Figures 8-13). An inflation of p-values compared to permuted phenotypes was observed only

182 within the tail end of highly significant SNPs. This is a characteristic not of population
183 stratification but of a trait being polygenic, i.e. with many truly causal SNPs each explaining
184 only a small proportion of variance. This distribution is common among human GWAS QQ
185 plots and suggests larger studies of DR will yield additional causal SNPs, albeit with smaller
186 effect sizes.

187 Usually population stratification is addressed by correcting for ancestry informative
188 principal components. These principal components are based on SNP correlations across the
189 genome, and have been shown to accurately capture population structure[23]. However their
190 construction proved difficult in our total sample due to much higher missingness than is
191 typical in GWAS data from genotyping chips. As such we performed a sensitivity analysis in
192 a smaller sample with near complete sequencing (n=178) to test the effect of our genome-
193 wide significant SNPs after correcting for principal components. No large attenuation of
194 effect was observed, with half of the genome-wide significant SNPs showing an increased
195 effect size when the first five principal components were included as covariates
196 (Supplementary Table 1). Predictably we observed higher p-values in the sensitivity test due
197 to the much smaller sample size. The partial availability of GPS data for individual's
198 household allowed for comparison of geographic proximity to genetic similarity (n=34). We
199 did not observe clear genetic clustering overlapping with geographic (see Figure 3), though a
200 pairwise comparison of genetic distance based on coordinates of first 2 principal components
201 and geographic position did show a weak association between the two ($R^2=1.4\%$, $p<0.005$).

202 Another potential confounder within GWAS is relatedness between samples.
203 Traditional measures of human relatedness were not appropriate for the analysis of pathogen
204 genomics data. We performed a sensitivity test to remove samples closely linked within
205 phylogenetic clusters (N=6). The results did not differ greatly, suggesting our top findings

206 were not driven by population stratification or cryptic relatedness (see Supplementary Table
207 1).

208 **Discussion**

209 In this study, we performed a proof of concept analysis that shows how a GWAS
210 approach can identify many known variants and replicable novel associations using HIV
211 WGS. We identified five variants at loci which corresponded with amino acid changes
212 previously associated with DR. While not all previously known DR variants were identified
213 at genome-wide significance in our analyses, we observed an excess of nominally significant
214 associations at these loci ($p < 0.001$, Supplementary Table 2). This is reminiscent of the
215 polygenicity observed in human GWAS. Often an excess of sub-genome-wide significant
216 variants was identified prior to identifying those specific SNPs truly associated with a
217 trait[24]. We can expect many of those previously known variants to become genome-wide
218 significant once sample sizes increase.

219 As well as validating known variants, our results highlight two ways in which GWAS
220 can identify potential novel variants. The first is by identifying variants of smaller or indirect
221 effects, such as via epistasis. We identified two nonsynonymous variants changing the RT
222 amino acid 68 from a serine to asparagines or glycine. Both associations remained after
223 correction for other treatments and potential confounders and the amino acid changes were
224 replicated in independent samples. The 68 glycine variant has been described previously as
225 correlating with drug resistance variant at position 65[21]. This change does not confer drug
226 resistance itself but rather compensates for the reduced fitness from a change at position
227 65[22]. In agreement with this we observed significant interactions between the changes at
228 position 68 and both 65 and 106 (Supplementary Table 3).

229 The second benefit of a GWAS approach was the ability to identify novel associations
230 outside of candidate regions of the genome. Here we observed a novel associated SNP
231 outside of the RT region of the Pol gene traditionally assumed to contain all genetic variants
232 that provide resistance to NNRTIs. This association with failure on tenofovir (an NNRTI)
233 was instead within amino acid 91 of the matrix protein of the Gag polyprotein. The effect
234 remained after correction for effects of other drugs, population stratification and relatedness
235 (Supplementary Table 4). This variant leads to a change in amino acid from the reference
236 arginine to glycine, an uncommon change though the region is highly polymorphic. In
237 comparison to the other variants (mean odds ratio of 3.70, range 1.72-11.91), the effect size
238 was slightly smaller at 1.78 suggesting why it previously may have been unobserved.

239 While the current results validate the applicability of GWAS to the HIV genome,
240 there are some limitations. As previously mentioned, not all known DR variants were
241 identified at genome-wide significance, though given many were nominally significant, this
242 is likely to reflect small sample size. Related to this is a bias in which types of variants were
243 more likely to be identified in our study design. These would have related to two groups of
244 variants. First, we would have had greater power to detect drug resistance variants that also
245 reduce viral fitness, meaning they would only exist at high frequencies when directly under
246 selection from treatment. Second, our study design would favour identifying variants that had
247 effects specific to one drug rather than a class of drugs, due to most samples having been
248 exposed to at least one drug from each class. This was a result of the now widespread usage
249 of ART by infected individuals and subsequent focus of sequencing efforts on treatment
250 resistance. Lastly, we note that unlike bacterial GWAS[7], we did not observe dramatic
251 genome-wide inflation in test statistics. Our comparison of lambda values using permuted
252 and unadjusted p-values suggested that Bonferroni adjustment for multiple corrections is
253 likely over conservative, while permutation adjustment may not correct for all inflation.

254 However, analysis of principle components suggested the genome-wide associations were not
255 confounded by geographic and genetic population structure.

256 Overall, our results provide a clear proof of concept on the use of GWAS within HIV
257 and other viruses whole genome sequence data. The smaller genome size, compared to
258 humans, means that substantially smaller samples were needed to identify associated variants.
259 Power is also greater because sequencing allows one to test the association with the causal
260 variant, rather than the proxy SNPs often used in human GWAS to capture several nearby
261 correlated SNPs. With a larger percentage of the genome transcribed there should also be a
262 larger proportion of functionally relevant variants. Additionally, viruses can themselves be
263 used as model organisms and can be genetically modified, allowing for functional validation
264 of identified variants in a way that cannot be performed in humans. However, these benefits
265 of performing GWAS within viruses should not ignore the valuable lessons from human
266 genomics, especially the need to quickly establish large sample sizes through internationally
267 collaborative research (see [25, 26]). A focus on setting up standardised quality control
268 pipelines, making GWAS results publically available in the form of SNP summary statistics,
269 and pooling samples into mega-analyses (rather than meta-analysing separate studies) should
270 be the aim of those groups generating HIV and other virus genomes.

271

272 **Methods**

273 *Sample description*

274 The study sampled 319 HIV-infected adults and 24 children on ART with virological
275 failure in the Hlabisa HIV Treatment and Care Programme in South Africa for which a whole
276 genome of HIV-1 was produced. The inclusion criteria were: ART regimen for at least 12

277 months followed by virological failure, defined as one viral load >1000 copies/ml. Exclusion
278 criteria were: prior use of nucleoside reverse transcriptase inhibitor (NRTI) monotherapy or
279 dual therapy (not including regimens for the prevention of mother-to-child transmission
280 (pMTCT)). All individuals were seen by a physician, who performed a clinical evaluation
281 and obtained written informed consent for the study. A 5 ml EDTA whole blood sample for
282 HIV DR genotyping was collected during the clinical evaluation. Basic clinical and
283 demographic data, including GPS data on household location, were collected on a clinical
284 form and clinical and treatment information was compared with the records in the Africa
285 Centre's ART Evaluation and Monitoring System (ARTEMIS), an operational database
286 holding treatment and laboratory monitoring information from the national ART programme
287 in South Africa. The clinical information was entered in anonymised form into a relational
288 sequence database, the SATuRN REGA database[27]. Further details of the study have been
289 described previously[28, 29].

290 *Ethics Statement*

291 The study was approved by the Biomedical Research Ethics Committee of the
292 University of KwaZulu-Natal (ref. BF052/10) and the Health Research Committee of the
293 KwaZulu-Natal Department of Health (ref. HRKM 176/10). South African legal guidelines
294 define a person able to give informed from consent from age 17. Written informed consent
295 was obtained from all the study participants and their parent or legal guardian in the case of
296 paediatric patients (≤ 16 years).

297 *Drug exposure data*

298 The median duration of ART among patients in this cohort was 42 months (IQR 32–
299 53). The most common first line ART regimens were: tenofovir/stavudine/zidovudine
300 +Lamivudine +efavirenz/nevirapine. The most common second line ART regimen were:

301 Lopinavir (+ Ritonavir), Lamivudine, zidovudine/tenofovir. The median duration of
302 antiretroviral failure was 27 months (IQR 17– 40 months). Details on drug exposure data and
303 DR results have been described previously [28, 29]. Drug exposure was defined by exposure
304 at any time point prior to sequencing. Table 1 provides a basic description of the
305 characteristics of the 343 individuals with viral WGS data included in the analysis.

306 ***RNA extraction, PCR amplification and whole genome Sequencing***

307 RNA was extracted from samples using the manual QIAamp Viral RNA Mini Kit
308 (Qiagen). The near complete HIV-1 genome was amplified by a previously described RT-
309 PCR strategy with primers modified to be more subtype C specific (Danaviah et al. CROI
310 2015; Abstract). The amplification involved the production of four overlapping genetic
311 fragments of lengths of 1.9kb, 3.6kb, 3.0kb and 3.5kb. This included all nine open reading
312 frames and partial regions of the 5'- and 3'-LTR. The DNA concentration of individual
313 amplicons was quantified using the Qubit ssDNA HS Assay Kit (Thermo Fischer Scientific-
314 Life Technologies). Pooled amplicons were prepared for sequencing using the Nextera XT
315 DNA Sample Preparation kit (Illumina) and the Nextera XT DNA Sample Preparation Index
316 Kit (Illumina), following the manufacture's protocol. The runs comprised pools of 96
317 samples that included three controls (one negative sample, one inter-run and one intra-run
318 control). All processes to generate WGS were undertaken locally at the Africa Centre
319 laboratory, Nelson R Mandela Medical School, University of KwaZulu-Natal, South Africa.

320 ***Bioinformatics pipeline: Whole genome quality control, assembly and phylogenetic*** 321 ***analysis***

322 Fastq quality control was performed using FASTQC(0.11.3) and QUASR(3.1)
323 software applications. Reads of less than 100bp in length and a quality score lower than 30
324 were excluded. In addition, the reads were trimmed up to 10bp from 5' and 30bp at the 3' to

325 exclude poor quality sequence at the beginning and end of reads. We noticed that the second
326 pair read of the Illumina Nextera XT was of lower quality and that excluding the last 30bp
327 increased quality score to > 33 . We imposed these exclusion criteria in order to decrease the
328 probability of ambiguous read mapping, which occurs when shorter reads of lower accuracy
329 are included in assemblies [30]. Following these quality control steps, we mapped reads
330 against a subtype C reference sequence (AF411967) with five assembly iterations using
331 Geneious 8 (<http://www.geneious.com>)[31]. After assembly, we exported the data as BAM
332 files and exported contigs as FASTA files.

333 In order to determine if there was clustering of sequences (i.e. sequences that were
334 very similar with low genetic diversity), we aligned all of the whole genomes with a
335 reference dataset for HIV-1 subtype C. The tree was constructed with HKY+Gamma site rate
336 variation in a MPI version of RaxML. Reliability of internal nodes was evaluated by 100
337 bootstrap replicates. Phylogenies were analysed using Phylotype software application[32] in
338 order to detect any clustering of sequences with high bootstrap values ($>90\%$) and low
339 sequence diversity ($<3\%$). This was performed to identify pairs of closely related HIV
340 sequences that might confound the analysis and test the sensitivity of the results to their
341 inclusion.

342 ***Variant calling and GWAS software adaptation***

343 The processing of WGS data to the performing of GWAS is outlined in Figure 1, with
344 comparison to human GWAS steps. BAM files were converted to VCF format variant calls
345 individually for each sequence in PILON [33]. A threshold of a depth of 50 reads per base
346 was used for a variant to be called.

347 As GWAS software was originally designed for diploid organisms (i.e. those with two
348 chromosomes and so two copies of any given loci), each sample can be called either as

349 homozygous for an allele (e.g. AA or TT) or as heterozygous (e.g. AT). While heterozygosity
350 is incorrect in the sense that HIV is haploid, it captures an important reality of viral infection:
351 genetic differences within the host's viral population. We wanted to retain the feature of
352 diploidy to account for samples with diversity at a given DR loci. We expected heterozygous
353 samples to have an intermediate effect size compared to samples where the DR variant was
354 either entirely non-existent or fixed. The downside of this approach was that given numerous
355 sequence reads for each loci, some variation is expected due to sequencing errors. To account
356 for this, we allowed for diploid calling in the following manner. If the reference allele
357 frequency was present in >85% of reads at a loci, the loci was called as homozygous for the
358 reference allele. A heterozygous call with one copy of the reference variant and one of the
359 non-reference variant was made if the reference allele frequency was between 85% and 15%
360 of reads. Finally, a homozygous non-reference call was made if the reference allele frequency
361 was found in less than 15% of reads. While these cut-offs are simply defaults of the software,
362 this worked as a crude calling approach for whether an individual sample's HIV population
363 was fixed or mixed for any given loci.

364 VCFs were then merged in GATK[34], then the combined VCF read into
365 PLINK1.09[35] for GWAS analysis. Prior to analysis, several QC steps were performed.
366 First, where multiple alleles occurred at the same loci, the reference variant and the most
367 common non-reference variant were used to make the loci bi-allelic. Second, a minor allele
368 frequency of greater than 1% was required for all variants. Lastly, we did not implement a
369 restriction on missingness of data. In human GWAS, high missingness for a SNP or
370 individual may reflect poor quality genotyping. However, in HIV WGS sequencing quality is
371 not homogenous across the genome (see Supplementary Figure 2). As we had restricted
372 analysis to calling variants at loci with a depth of 50 or greater, higher missingness was
373 expected. Missingness for SNPs significantly associated with DR is reported in Table 2.

374 ***Statistical analysis***

375 A logistic regression was performed in PLINK1.09[35] with drug exposure as the
376 binary outcome and each SNP as a predictor with an additive effect. All samples exposed to a
377 given drug were compared to all that were not. To determine genome-wide significance we
378 performed 10,000,000 permutations within PLINK1.09 both on a single SNP and genome-
379 wide level using the --mperm command. This was performed to account for correlation
380 between nearby SNPs which would have made Bonferroni correction for the raw number of
381 statistical tests overly conservative. Given the smaller number of variants compared to a
382 human GWAS, permutation using 10,000,000 for the empirical p-values was computationally
383 feasible. As the negative correlations in the prescribing of these drugs existed, associations
384 with the same SNP were seen in multiple analyses. However it was possible to identify when
385 exposure was associated with the non-reference sequence (i.e. odds ratio>1) and so,
386 presumably, which association identified the true drug resistant variant. Principal components
387 were generated in GCTA[36].

388 ***Replication***

389 Genome-wide significant SNPs within the Pol region were able to be taken forward
390 for replication in a publically available independent sample. This was the Stanford University
391 HIV Drug Resistance Database[19, 20], where information on amino acid frequencies were
392 available for sequences exposed to different drugs. This analysis was restricted to the 13,676
393 subtype C sequences. Additional analyses also made use of a subset of all publically available
394 subtype C WGS (n=505) from the Los Alamos HIV Sequence Database
395 (<http://www.hiv.lanl.gov/>). This was done to ensure our variant frequencies were in
396 agreement with those observed elsewhere.

397 ***Data access***

398 Summary statistics for all SNPs of each GWAS are available online
399 (https://figshare.com/articles/PLOSONE_DR_GWAS_HIV/3569766). Access to the full
400 genomes of HIV can be done by application of a proposal to PANGEA_HIV
401 (<http://www.pangea-hiv.org/>).

402

403

Funding Disclosure

404 Research supported by a South African MRC Flagship grant (MRC-RFA-UFSP-01-
405 2013/UKZN HIVEPI) and the Bill and Melinda Gates Foundation (BMGF) PANGEA_HIV
406 grant. Funding for the Africa Centre's Demographic Surveillance Information System and
407 Population-Based HIV Survey was received from the Wellcome Trust, UK (grant
408 082384/Z/07/Z). The funders had no role in study design, data collection and analysis,
409 decision to publish, or preparation of the manuscript.

410

411

Acknowledgements

412 We thank all the patients for their continued support, all Africa Centre staff who contribute to
413 maintaining ACDIS and the laboratory.

414

415

Competing Interests

416 The authors have declared that no competing interests exist.

417

418

References

- 419 1. Bush WS, Moore JH. Chapter 11: Genome-wide association studies. *PLoS Comput Biol*.
420 2012;8(12):e1002822. doi: 10.1371/journal.pcbi.1002822. PubMed PMID: 23300413; PubMed
421 Central PMCID: PMC3531285.
- 422 2. Jostins L, Ripke S, Weersma RK, Duerr RH, McGovern DP, Hui KY, et al. Host-microbe
423 interactions have shaped the genetic architecture of inflammatory bowel disease. *Nature*.
424 2012;491(7422):119-24. doi: 10.1038/nature11582. PubMed PMID: 23128233; PubMed Central
425 PMCID: PMC3491803.
- 426 3. Wood AR, Esko T, Yang J, Vedantam S, Pers TH, Gustafsson S, et al. Defining the role of
427 common variation in the genomic and biological architecture of adult human height. *Nat Genet*.
428 2014;46(11):1173-86. doi: 10.1038/ng.3097. PubMed PMID: 25282103; PubMed Central PMCID:
429 PMC4250049.
- 430 4. Global Lipids Genetics C, Willer CJ, Schmidt EM, Sengupta S, Peloso GM, Gustafsson S, et al.
431 Discovery and refinement of loci associated with lipid levels. *Nat Genet*. 2013;45(11):1274-83. doi:
432 10.1038/ng.2797. PubMed PMID: 24097068; PubMed Central PMCID: PMC3838666.
- 433 5. Pillay D, Herbeck J, Cohen MS, de Oliveira T, Fraser C, Ratmann O, et al. PANGEA-HIV:
434 phylogenetics for generalised epidemics in Africa. *Lancet Infect Dis*. 2015;15(3):259-61. doi:
435 10.1016/S1473-3099(15)70036-8. PubMed PMID: 25749217.
- 436 6. Malaria Genomic Epidemiology N. A global network for investigating the genomic
437 epidemiology of malaria. *Nature*. 2008;456(7223):732-7. doi: 10.1038/nature07632. PubMed PMID:
438 19079050; PubMed Central PMCID: PMC3758999.
- 439 7. Read TD, Massey RC. Characterizing the genetic basis of bacterial phenotypes using genome-
440 wide association studies: a new direction for bacteriology. *Genome Med*. 2014;6(11):109. doi:
441 10.1186/s13073-014-0109-z. PubMed PMID: 25593593; PubMed Central PMCID: PMC4295408.
- 442 8. Miotto O, Amato R, Ashley EA, MacInnis B, Almagro-Garcia J, Amaratunga C, et al. Genetic
443 architecture of artemisinin-resistant *Plasmodium falciparum*. *Nat Genet*. 2015;47(3):226-34. doi:
444 10.1038/ng.3189. PubMed PMID: 25599401.
- 445 9. Farhat MR, Shapiro BJ, Kieser KJ, Sultana R, Jacobson KR, Victor TC, et al. Genomic analysis
446 identifies targets of convergent positive selection in drug-resistant *Mycobacterium tuberculosis*. *Nat*
447 *Genet*. 2013;45(10):1183-9. doi: 10.1038/ng.2747. PubMed PMID: 23995135; PubMed Central
448 PMCID: PMC3887553.
- 449 10. Alam MT, Petri RA, 3rd, Crispell EK, Thornton TA, Conneely KN, Jiang Y, et al. Dissecting
450 vancomycin-intermediate resistance in *Staphylococcus aureus* using genome-wide association.
451 *Genome Biol Evol*. 2014;6(5):1174-85. doi: 10.1093/gbe/evu092. PubMed PMID: 24787619; PubMed
452 Central PMCID: PMC4040999.
- 453 11. Chewapreecha C, Marttinen P, Croucher NJ, Salter SJ, Harris SR, Mather AE, et al.
454 Comprehensive identification of single nucleotide polymorphisms associated with beta-lactam
455 resistance within pneumococcal mosaic genes. *PLoS Genet*. 2014;10(8):e1004547. doi:
456 10.1371/journal.pgen.1004547. PubMed PMID: 25101644; PubMed Central PMCID: PMC4125147.
- 457 12. Bartha I, Carlson JM, Brumme CJ, McLaren PJ, Brumme ZL, John M, et al. A genome-to-
458 genome analysis of associations between human genetic variation, HIV-1 sequence diversity, and
459 viral control. *Elife*. 2013;2:e01123. doi: 10.7554/eLife.01123. PubMed PMID: 24171102; PubMed
460 Central PMCID: PMC3807812.
- 461 13. Didelot X, Maiden MC. Impact of recombination on bacterial evolution. *Trends Microbiol*.
462 2010;18(7):315-22. doi: 10.1016/j.tim.2010.04.002. PubMed PMID: 20452218; PubMed Central
463 PMCID: PMC3985120.
- 464 14. UNAIDS. Together we will end AIDS. Geneva: UNAIDS, 2012.
- 465 15. Bor J, Herbst AJ, Newell ML, Barnighausen T. Increases in adult life expectancy in rural South
466 Africa: valuing the scale-up of HIV treatment. *Science*. 2013;339(6122):961-5. doi:
467 10.1126/science.1230413. PubMed PMID: 23430655; PubMed Central PMCID: PMC3860268.

- 468 16. Tanser F, Barnighausen T, Grapsa E, Zaidi J, Newell ML. High coverage of ART associated with
469 decline in risk of HIV acquisition in rural KwaZulu-Natal, South Africa. *Science*. 2013;339(6122):966-
470 71. doi: 10.1126/science.1228160. PubMed PMID: 23430656; PubMed Central PMCID: PMC4255272.
- 471 17. Cohen MS, Chen YQ, McCauley M, Gamble T, Hosseinipour MC, Kumarasamy N, et al.
472 Prevention of HIV-1 infection with early antiretroviral therapy. *N Engl J Med*. 2011;365(6):493-505.
473 doi: 10.1056/NEJMoa1105243. PubMed PMID: 21767103; PubMed Central PMCID: PMC3200068.
- 474 18. Wensing AM, Calvez V, Gunthard HF, Johnson VA, Paredes R, Pillay D, et al. 2014 Update of
475 the drug resistance mutations in HIV-1. *Top Antivir Med*. 2014;22(3):642-50. PubMed PMID:
476 25101529; PubMed Central PMCID: PMC4392881.
- 477 19. Rhee SY, Gonzales MJ, Kantor R, Betts BJ, Ravela J, Shafer RW. Human immunodeficiency
478 virus reverse transcriptase and protease sequence database. *Nucleic Acids Res*. 2003;31(1):298-303.
479 PubMed PMID: 12520007; PubMed Central PMCID: PMC165547.
- 480 20. Shafer RW. Rationale and uses of a public HIV drug-resistance database. *J Infect Dis*.
481 2006;194 Suppl 1:S51-8. doi: 10.1086/505356. PubMed PMID: 16921473; PubMed Central PMCID:
482 PMC2614864.
- 483 21. Margot NA, Waters JM, Miller MD. In vitro human immunodeficiency virus type 1 resistance
484 selections with combinations of tenofovir and emtricitabine or abacavir and lamivudine. *Antimicrob*
485 *Agents Chemother*. 2006;50(12):4087-95. doi: 10.1128/AAC.00816-06. PubMed PMID: 16982781;
486 PubMed Central PMCID: PMC1693985.
- 487 22. Svarovskaia ES, Feng JY, Margot NA, Myrick F, Goodman D, Ly JK, et al. The A62V and S68G
488 mutations in HIV-1 reverse transcriptase partially restore the replication defect associated with the
489 K65R mutation. *J Acquir Immune Defic Syndr*. 2008;48(4):428-36. doi:
490 10.1097/QAI.0b013e31817bbe93. PubMed PMID: 18614922.
- 491 23. Novembre J, Johnson T, Bryc K, Kutalik Z, Boyko AR, Auton A, et al. Genes mirror geography
492 within Europe. *Nature*. 2008;456(7218):98-101. doi: 10.1038/nature07331. PubMed PMID:
493 18758442; PubMed Central PMCID: PMC2735096.
- 494 24. Purcell SM, Wray NR, Stone JL, Visscher PM, O'Donovan MC, Sullivan PF, et al. Common
495 polygenic variation contributes to risk of schizophrenia and bipolar disorder. *Nature*.
496 2009;460(7256):748-52. Epub 2009/07/03. doi: nature08185 [pii]
10.1038/nature08185. PubMed PMID: 19571811.
- 497 25. Gratten J, Wray NR, Keller MC, Visscher PM. Large-scale genomics unveils the genetic
498 architecture of psychiatric disorders. *Nat Neurosci*. 2014;17(6):782-90. doi: 10.1038/nn.3708.
499 PubMed PMID: 24866044; PubMed Central PMCID: PMC4112149.
- 500 26. Ahlqvist E, van Zuydam NR, Groop LC, McCarthy MI. The genetics of diabetic complications.
501 *Nat Rev Nephrol*. 2015;11(5):277-87. doi: 10.1038/nrneph.2015.37. PubMed PMID: 25825086.
- 502 27. Manasa J, Lessells R, Rossouw T, Naidu K, Van Vuuren C, Goedhals D, et al. Southern African
503 Treatment Resistance Network (SATuRN) RegaDB HIV drug resistance and clinical management
504 database: supporting patient management, surveillance and research in southern Africa. *Database*
505 (Oxford). 2014;2014:bat082. doi: 10.1093/database/bat082. PubMed PMID: 24504151.
- 506 28. Manasa J, Lessells RJ, Skingsley A, Naidu KK, Newell ML, McGrath N, et al. High-levels of
507 acquired drug resistance in adult patients failing first-line antiretroviral therapy in a rural HIV
508 treatment programme in KwaZulu-Natal, South Africa. *PLoS One*. 2013;8(8):e72152. doi:
509 10.1371/journal.pone.0072152. PubMed PMID: 23991055; PubMed Central PMCID: PMC3749184.
- 510 29. Pillay S, Bland RM, Lessells RJ, Manasa J, de Oliveira T, Danaviah S. Drug resistance in
511 children at virological failure in a rural KwaZulu-Natal, South Africa, cohort. *AIDS Res Ther*.
512 2014;11(1):3. doi: 10.1186/1742-6405-11-3. PubMed PMID: 24444369; PubMed Central PMCID:
513 PMC3922737.
- 514 30. Watson SJ, Welkers MR, Depledge DP, Coulter E, Breuer JM, de Jong MD, et al. Viral
515 population analysis and minority-variant detection using short read next-generation sequencing.
516 *Philos Trans R Soc Lond B Biol Sci*. 2013;368(1614):20120205. doi: 10.1098/rstb.2012.0205. PubMed
517 PMID: 23382427; PubMed Central PMCID: PMC3678329.
- 518

- 519 31. Kearse M, Moir R, Wilson A, Stones-Havas S, Cheung M, Sturrock S, et al. Geneious Basic: an
520 integrated and extendable desktop software platform for the organization and analysis of sequence
521 data. *Bioinformatics*. 2012;28(12):1647-9. doi: 10.1093/bioinformatics/bts199. PubMed PMID:
522 22543367; PubMed Central PMCID: PMC3371832.
- 523 32. Chevenet F, Jung M, Peeters M, de Oliveira T, Gascuel O. Searching for virus phlotypes.
524 *Bioinformatics*. 2013;29(5):561-70. doi: 10.1093/bioinformatics/btt010. PubMed PMID: 23329414;
525 PubMed Central PMCID: PMC3582263.
- 526 33. Walker BJ, Abeel T, Shea T, Priest M, Abouelliel A, Sakthikumar S, et al. Pilon: an integrated
527 tool for comprehensive microbial variant detection and genome assembly improvement. *PLoS One*.
528 2014;9(11):e112963. doi: 10.1371/journal.pone.0112963. PubMed PMID: 25409509; PubMed
529 Central PMCID: PMC4237348.
- 530 34. McKenna A, Hanna M, Banks E, Sivachenko A, Cibulskis K, Kernytisky A, et al. The Genome
531 Analysis Toolkit: a MapReduce framework for analyzing next-generation DNA sequencing data.
532 *Genome Res*. 2010;20(9):1297-303. doi: 10.1101/gr.107524.110. PubMed PMID: 20644199; PubMed
533 Central PMCID: PMC2928508.
- 534 35. Chang CC, Chow CC, Tellier LC, Vattikuti S, Purcell SM, Lee JJ. Second-generation PLINK: rising
535 to the challenge of larger and richer datasets. *Gigascience*. 2015;4:7. doi: 10.1186/s13742-015-0047-
536 8. PubMed PMID: 25722852; PubMed Central PMCID: PMC4342193.
- 537 36. Yang J, Lee SH, Goddard ME, Visscher PM. GCTA: a tool for genome-wide complex trait
538 analysis. *Am J Hum Genet*. 2011;88(1):76-82. doi: 10.1016/j.ajhg.2010.11.011. PubMed PMID:
539 21167468; PubMed Central PMCID: PMC3014363.

540

541

542

Tables and Figures

Table 1: Number of WGS treated with each drug, and correlations between drugs within samples

| Drug | Treated | Untreated | Correlation with: | | | | | |
|-------------------|---------|-----------|-------------------|-----------|-----------|-----------|------------|-----------|
| | | | Zidovudine | Stavudine | Tenofovir | Efavirenz | Nevirapine | Lopinavir |
| Zidovudine | 32 | 311 | 1 | - | - | - | - | - |
| Stavudine | 291 | 52 | -0.058 | 1 | - | - | - | - |
| Tenofovir | 101 | 242 | -0.117 | -0.507 | 1 | - | - | - |
| Efavirenz | 259 | 84 | 0.011 | -0.023 | 0.128 | 1 | - | - |
| Nevirapine | 113 | 230 | -0.017 | 0.127 | -0.057 | -0.623 | 1 | - |
| Lopinavir | 26 | 317 | 0.213 | -0.033 | 0.053 | -0.151 | -0.115 | 1 |

Table 2: Results for genome-wide significant SNPs and their corresponding amino acid positions. Note that the effect of SNP 2739A is protective against stavudine resistance (i.e. odds ratio [OR] <1) and the association is actually with tenofovir, that has a negatively correlated prescription regime. Ref.=Reference; BP= base position; A1 = effect allele; Cis=proximal to known DR variant; Conv.=convergent, i.e. known DR variant for another drug; OR=Odds ratio; SE=standard error.

| Drug | SNP | Missing -ness | A1 | Ref. | Gene | Amino acid N | Ref. Amino Acid | A1 Amino Acid | Known | OR | SE | Unadjust ed p- value | Permutatio n adjusted p-value |
|------------|-------|------------------|----|------|----------|-----------------|--------------------|------------------|-------|------|------|----------------------------|-------------------------------------|
| Nevirapine | 3078G | 14% | G | A | RT | 181 | Y | C | Yes | 5.20 | 0.26 | 4.77E-10 | 1.00E-07 |
| Stavudine | 2739A | 14% | A | G | RT | 68 | S | N | Cis | 0.08 | 0.54 | 5.38E-06 | 0.0081 |
| Tenofovir | 1063A | 18% | G | A | MA (p17) | 91 | R | G | No | 1.79 | 0.14 | 2.42E-05 | 0.016 |
| | 2730G | 13% | G | A | RT | 65 | K | R | Yes | 6.44 | 0.24 | 1.67E-14 | 1.00E-07 |
| | 2738G | 14% | G | A | RT | 68 | S | G | Cis | 2.89 | 0.24 | 1.45E-05 | 0.0088 |
| | 2852A | 14% | A | G | RT | 106 | V | M | Conv. | 1.72 | 0.14 | 6.19E-05 | 0.047 |
| | 2880T | 13% | T | A | RT | 115 | Y | F | Conv. | 5.77 | 0.41 | 1.80E-05 | 0.011 |
| Zidovudine | 2745G | 16% | G | A | RT | 70 | K | R | Yes | 3.11 | 0.22 | 2.94E-07 | 0.0006 |

Figure 1: Analysis pipeline for HIV whole genome sequence (WGS) genome-wide association study (GWAS) compared to a human study using a SNP chip. Step 1) Diploidy defined for both human and pathogen, to reflect ‘real’ heterozygosity and heterozygosity from within host viral diversity. 2) While missingness and Hardy-Weinberg Equilibrium are used to assess genotyping quality in human GWAS, in viral GWAS we used depth of sequencing to assess variant calls. As such, higher calling confidence is associated with higher missingness in viral SNPs, while the reverse is true in humans. Low minor allele frequency (MAF) is always used to remove variants that have low power to detect effects and may reflect errors. 3&4) Correction for ancestry and relatedness are key to human GWAS, however due to both more homogenous sampling and difficulty in applying conventional corrections in human data to viral, this was done as a sensitivity test in a smaller sample for top SNPs in HIV GWAS.

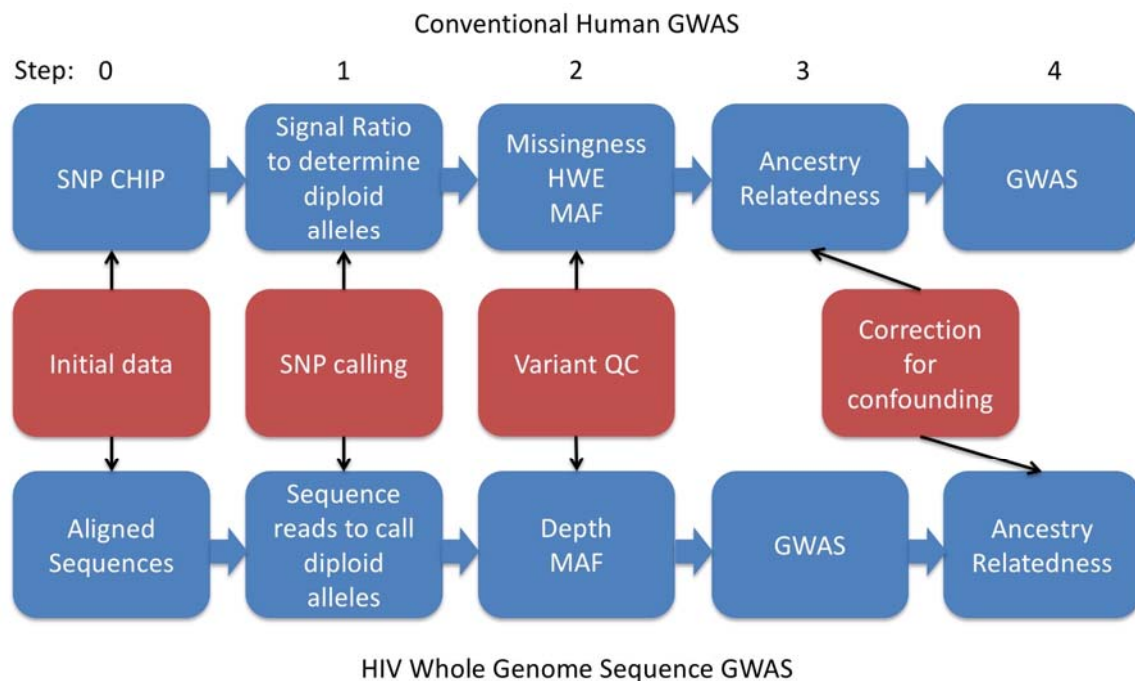


Figure 2: Manhattan plot comparing HIV sequences that were exposed to tenofovir to those that were not. The reference line at $p=7E-5$ is the line for permutation adjusted genome wide significance. Dashed grey lines on genomic locations refer to borders of genes (black dashed refer to GAG, Pol and ENV). Each marker is a SNP, weighted by its $-\log(p\text{-value})$ to highlight the most significant SNPs.

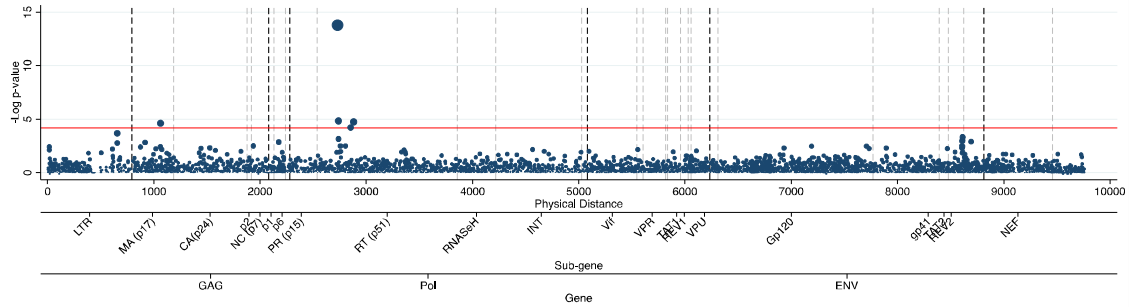


Figure 3: Plot of standardised values for the ancestry informative principle components 1 & 2 (red) and latitude & longitude (gold) for HIV sequences, with values for each sequence linked by a line. No correlation between geographic position and genetic position was observed.

