

1 **Classification:**

2 Major: Biological Sciences

3 Minor: Genetics

4

5 **HIV Peptidome-Wide Association Study Reveals Patient-Specific Epitope**

6 **Repertoires Associated with HIV Control**

7

8 **Short Title: Epitope Repertoires Associated with HIV Control**

9

10

11 Jatin Arora¹, Paul J. McLaren^{2,3}, Nimisha Chaturvedi^{4,5}, Mary Carrington^{6,7},

12

Jacques Fellay^{4,5} & Tobias L. Lenz^{1,*}

13

14 **Affiliation Affiliations:**

15 ¹ Research Group for Evolutionary Immunogenomics, Max Planck Institute for

16 Evolutionary Biology, 24306 Plön, Germany. ² JC Wilt Infectious Diseases Research

17 Center, National HIV and Retrovirology Laboratory, Public Health Agency of Canada,

18 R3E 0W3, Winnipeg, Canada. ³ Department of Medical Microbiology and Infectious

19 Diseases, University of Manitoba, R3E 0J9, Winnipeg Canada. ⁴ Global Health Institute,

20 School of Life Sciences, École Polytechnique Fédérale de Lausanne, 1015 Lausanne,

21 Switzerland. ⁵ Swiss Institute of Bioinformatics, 1015 Lausanne, Switzerland. ⁶ Cancer and

22 Inflammation Program, Leidos Biomedical Research, Frederick National Laboratory,

23 Frederick, MD 21702, USA. ⁷ Ragon Institute of Massachusetts General Hospital,

24 Massachusetts Institute of Technology and Harvard University, Cambridge, MA 02139-

25 3583, USA

26

27 **Corresponding Author:**

28 Tobias L. Lenz

29 Max Planck Institute for Evolutionary Biology

30 August-Thienemann-Str. 2

31 Plön 24306, Germany

32 Phone: +49 4522 763 228

33 Email: lenz@post.harvard.edu

34

35 **Keywords:**

36 Human Leukocyte Antigen (HLA) | HIV | Evolution | Adaptive Immunity

37

38 **Abstract:**

39 Genetic variation in the peptide-binding groove of the highly polymorphic human
40 leukocyte antigen (HLA) class I molecules has repeatedly been associated with HIV-1
41 control and progression to AIDS, accounting for up to 12% of the variation in HIV-1 set
42 point viral load (spVL). This suggests a key role in disease control for HLA presentation
43 of HIV-1 epitopes to cytotoxic T cells. However, a comprehensive understanding of the
44 relevant HLA-bound HIV epitopes is still elusive. Here we describe a peptidome-wide
45 association study (PepWAS) approach that integrates HLA genotypes and spVL data from
46 6,311 HIV-infected patients to interrogate the entire HIV-1 proteome (3,252 unique
47 peptides) for disease-relevant peptides. This PepWAS approach predicts a core set of
48 epitopes associated with spVL, including previously characterized epitopes but also several
49 novel disease-relevant peptides. More importantly, each patient presents only a small
50 subset of these predicted core epitopes through their individual HLA-A and -B variants.
51 Eventually, the individual differences in these patient-specific epitope repertoires account
52 for the variation in spVL that was previously associated with HLA genetic variation.
53 PepWAS thus enables a comprehensive functional interpretation of the robust but little
54 understood association between HLA and HIV-1 control, prioritizing a short list of disease-
55 associated epitopes for the development of targeted therapy.

56

57 **Significance Statement:**

58 Individual differences in HIV-1 control and progression to AIDS have been pinpointed to
59 genetic variation in the Human Leukocyte Antigen (HLA), coding for antigen-presenting
60 molecules. However, our understanding of the corresponding antigens is still incomplete.
61 Here we developed a new approach that combines HLA genotypes and viral load data of
62 HIV infected individuals to screen the entire HIV proteome for disease-relevant peptides.
63 Our PepWAS approach identified a limited manageable core set of peptides, accounting
64 for the entire variation in viral load previously associated with genetic variation in the
65 HLA. This core set of disease-relevant antigens thus provides a functional link between
66 HLA genetic variation and HIV-1 control, confirming several known antigens, but also
67 prioritizing novel antigens as new therapeutic targets.
68

69 \body
70 HLA class I proteins are thought to play a critical role in immune recognition of HIV-1 by
71 presenting endogenously processed viral peptides at the surface of infected cells to
72 cytotoxic T cells, in order to trigger destruction of the infected cells (1). Indeed, genetic
73 variation in the HLA region has repeatedly been identified as the major genetic determinant
74 of HIV-1 control in genome-wide association studies (2, 3). Most recently, McLaren *et al.*
75 (4) fine-mapped the entire HLA's association with HIV-1 control and disease progression
76 to five independent amino acid residues in the peptide binding groove of the HLA-B and
77 HLA-A molecules. These five residues alone accounted for 12.3% of the variation in viral
78 load, suggesting a major role for specific HLA-presented viral epitopes in HIV-1 control.
79 However, our understanding of the disease-relevant viral epitopes is still incomplete,
80 hampered by the economically hardly feasible challenge of employing a full-factorial
81 experimental assay to screen the entirety of the HIV-1 peptidome for binding by all relevant
82 HLA alleles. Therefore, we developed a novel computational analysis approach that
83 identifies and prioritizes disease-associated peptides based on individual HLA genotype
84 and disease phenotype information. Our approach uses established computational
85 algorithms to predict for each individual whether a given peptide is bound by the
86 individual's HLA variants, and then uses regression analysis on the disease phenotype
87 (here HIV set point viral load) to estimate whether the ability to bind the peptide is non-
88 randomly associated with the disease phenotype. This approach is analogous to a genetic
89 association study, except that it incorporates one additional layer by translating genetic
90 variation into functional variation (HLA variant-specific peptide binding). Importantly, this
91 approach does not simply define all peptides bound by a risk HLA variant as risk peptides.
92 *Instead, for each peptide it integrates the disease effect of all HLA variants that are able*
93 *to bind the peptide and thus estimates a peptide-specific association with disease.* Since
94 most peptides are bound by several HLA variants, integrating the effect of all binding HLA
95 variants is essential (**Fig. 1**). For instance, a peptide can have no association with disease,
96 even if it is bound by the highest risk variant, simply because it is also bound by several
97 other non-associated (or even protective) variants. Ultimately, our approach identifies a list
98 of peptides with varying associations to disease, which can directly inform therapy
99 development by prioritizing global as well as patient-specific candidate epitopes. As a

100 proof-of-concept, we analyze here a unique dataset of 6,311 individuals of European
101 ancestry with chronic HIV-1 infection (*SI Appendix, Table S1*). Screening the entire HIV-
102 1 peptidome for candidate epitopes, we identify a comprehensive list of peptides that
103 explain the well-established association between HLA genetic variation and HIV-1 control,
104 including several previously uncharacterized epitopes as novel candidates for targeted
105 therapy.

106

107 **Results and Discussion:**

108 Our analyses are based on a large dataset of HIV-infected individuals (4) that includes both
109 pre-treatment level of set point viral load (spVL) as a correlate of disease progression (5)
110 and imputed HLA genotypes (4-digit allele resolution). We focused on the two HLA loci
111 (HLA-B and HLA-A) reported to have independent associations with HIV-1 control and
112 disease progression (4). Potential HLA-bound peptides were identified using an established
113 computational algorithm that is based on empirical training data (6) and integrates several
114 complementary prediction methods in a consensus approach, outperforming comparable
115 algorithms (6, 7). Such algorithms have been used in a wide spectrum of HLA-related
116 studies ranging from vaccine design to cancer evolution and HIV disease genetics (8–10).
117 Without *a-priori* selection, we screened all possible 9mer HIV-1 peptides ($N = 3,252$) in a
118 sliding window across the entire HIV-1 M group subtype B reference proteome (11) against
119 all represented HLA-B and HLA-A alleles (344,712 HLA:peptide complexes), and
120 identified 214 and 173 distinct HIV-1 peptides predicted to be bound by one or more of the
121 represented HLA-B and HLA-A alleles, respectively.

122 In order to evaluate the significance of the predicted peptide repertoires, we interrogated
123 several layers of empirical evidence (see *SI Appendix, Supporting Text*). We observed an
124 enrichment for previously known HIV-1 epitopes (*SI Appendix, Fig. S1A*), a correlation
125 between an HLA-B allele's effect on viral load and the number of HIV-1 peptides it is
126 predicted to bind (*SI Appendix, Fig. S1B*), and detected previously reported viral escape
127 mutations (*SI Appendix, Fig. S1C*). Following these independent layers of evidence that
128 our analysis pipeline predicts disease-relevant binding of HLA to HIV-1 peptides, we
129 subsequently refer to the entire predicted set of HLA-bound peptides as *predicted epitopes*,
130 highlighting the point that not all of them have been experimentally validated.

131 Next, we tested whether the patient-specific repertoire of predicted HIV-1 epitopes, defined
132 by the number of peptides predicted to be bound by the specific HLA allele combination
133 of the patient, was associated with spVL. For this, we ran a linear regression across the
134 6,311 HIV-1 patients, with spVL as dependent variable and the patient-specific number of
135 bound peptides as predicting variable, together with other covariates (see methods). We
136 first focused on the effect of peptides bound by HLA-B, and used only the known CTL
137 epitopes from the Los Alamos HIV Molecular Immunology Database (12), of which 80
138 were represented among the 214 predicted HLA-B bound epitopes. The individual number
139 of these known CTL epitopes bound by patient-specific HLA-B variants accounted for only
140 1.8% of the individual variation in spVL (**Fig. 2**). In order to evaluate this association, we
141 then included all predicted HLA-bound HIV-1 epitopes (N = 214) in the analysis, including
142 the previously known CTL epitopes as well as any other HLA-B-bound peptide from the
143 HIV-1 proteome. Interestingly, the total number of all predicted HLA-B-bound epitopes
144 per patient accounted for 5.3% variation in spVL (**Fig. 2**), suggesting that the Los Alamos
145 CTL epitope dataset is not yet fully saturated with regard to disease-relevant peptides.
146 However, the accounted variation was still lower than the 11.4% variation associated with
147 genetic variation at HLA-B in previous genotype-based studies, suggesting that the total
148 predicted epitope repertoire still included peptides irrelevant for the association between
149 HLA and HIV-1. This is supported by a previous study, which showed that not all HLA-
150 bound peptides are epitopes targeted by CD8⁺ T-cells (13). We thus aimed to refine the
151 repertoire of predicted HLA-bound HIV-1 epitopes further to comprise only disease-
152 relevant epitopes. For this, we calculated the epitope-specific association with spVL by
153 running a separate linear regression for each predicted epitope and recording R² and β -
154 coefficient as measures of the epitope's effect on spVL. This is analogous to the approach
155 of a genome-wide association study (GWAS), where each genetic variant is tested for its
156 association with a given trait, except that here we focus on functional protein variation
157 (peptide binding by a patient's HLA molecules) rather than genetic variation. Following
158 this analogy, we term our approach *peptidome-wide association study* (PepWAS). Of 214
159 HIV-1 epitopes predicted to be bound by HLA-B, 132 accounted for nominal variation
160 (adjusted R² value > 0) in spVL, 74 of which were negatively and 58 positively associated
161 with spVL (β -coefficients ranging from -0.1 to 0.77; **SI Appendix, Table S2**). Importantly,

162 we do not require statistical significance at this point as this is a candidate screen and we
163 thus aim to minimize the number of false negatives. Subsequently, we designate the
164 nominally associated epitopes as *disease-associated* predicted epitopes, even though their
165 effects are not necessarily independent as they were tested with separate regression models.
166 An analogous investigation of peptide binding by HLA-A alleles revealed an additional 74
167 disease-associated epitopes (*SI Appendix, Table S3*).

168 Having refined the predicted HIV-1 epitope repertoire to only disease-associated predicted
169 epitopes, we then tested whether this subset accounted for a larger fraction of the variation
170 in spVL than the total predicted HIV-1 epitope repertoire. Indeed, the patients' ability to
171 bind a smaller or larger fraction of the HLA-B-specific disease-associated predicted
172 epitopes accounted for 11.4% of the variation in spVL (**Fig. 2**). Similarly, for HLA-A, the
173 total number of predicted HIV-1 epitopes bound by individual HLA-A genotypes
174 accounted for 0.3% of the variation, while disease-associated predicted epitopes accounted
175 for 1.4% of the variation in spVL. On average, a patient's HLA-B allele pair bound 16.2
176 ± 7 (SD) disease-associated predicted HIV-1 epitopes, while its HLA-A alleles bound
177 significantly less (6.6 ± 6.5 ; Paired Wilcoxon rank sum test, $P < 0.0001$; *SI Appendix, Fig.*
178 **S6**). This quantitative difference in peptide presentation might contribute to the stronger
179 spVL-association of HLA-B compared to HLA-A, as a larger number of presented peptides
180 should more likely lead to a more efficient CD8 T cell response, as has indeed been
181 observed for HLA-B compared to HLA-A (14). HLA-C-bound epitopes did not show any
182 significant association with spVL, mirroring the lack of independent genetic associations
183 for HLA-C in the latest GWAS (4). Predicted disease-associated epitopes of HLA-B and
184 HLA-A together accounted for 12.2% of the variation in HIV-1 viral load, approximately
185 corresponding to the 12.3% variation previously attributed to all independent genetic
186 associations in the entire HLA (**Fig. 2A**).

187 Interestingly, the *Env* protein showed the largest number of disease-associated predicted
188 epitopes, with both positive and negative effects. Among the disease-associated predicted
189 HLA-B-bound epitopes, *Env*-derived epitopes alone accounted for 6.4% of variation in
190 spVL, the highest among all HIV-1 proteins (**Fig. 3A**). In addition to already known *Env*-
191 derived CD8+ T-cell targeted epitopes associated with lower viral load and disease control
192 e.g. RIKQIINMW, HRLRDLILI (13), ERYLKDQQL (15), our analysis revealed

193 previously undescribed HLA-epitope complexes e.g. B*57:01-STQLFNSTW, -
194 NSTWFNSTW, or -RGWEALKYW showing strong associations with lower viral load
195 (**Fig. 3C**). The potential importance of the predicted *Env* epitopes is quite surprising, since
196 the high genetic variability of the *Env* protein across different HIV-1 isolates suggests that
197 the virus could readily evolve escape variants in this protein. However, a previous study
198 has already established that sequence conservation alone is not a reliable predictor of
199 protective epitopes, instead highlighting structural conservation as the more important
200 feature (13). More intriguingly, we found that the protective *Env* epitopes predicted through
201 our PepWAS approach are significantly enriched for residues that are associated with
202 broadly neutralizing antibodies (bNAbs; OR = 1.5, $P = 0.036$, *SI Appendix, Fig. S7*),
203 suggesting that they represent parts of the Env protein that can be efficiently targeted in
204 both antibody therapy as well as in HLA-mediated CTL response.

205 Notably, several of the represented HLA alleles were predicted to bind both negatively and
206 positively disease-associated epitopes (*SI Appendix, Tables S4 and S5*), i.e. epitopes
207 bound by the same HLA allele did not necessarily have the same effect on viral load. This
208 can be explained by the fact that a given epitope can be bound by several different HLA
209 alleles with very distinct disease association (see schematic in **Fig. 1**). This is also in
210 agreement with a previous study showing that viral control is mediated by specific
211 immunogenic epitopes which could be restricted by HLA alleles other than already known
212 ones (13).

213 HLA molecule variants are known to bind peptide repertoires with distinct anchor motifs,
214 based on the composition of their peptide-binding groove (16). This entailed the possibility
215 that our PepWAS approach is merely identifying distinct groups of peptides per HLA
216 variant, thus translating the known HLA variant-specific effect on viral load into peptide
217 group-specific effects. While still helpful in guiding epitope research, this would provide
218 only limited knowledge-gain compared to the HLA allele-specific associations known
219 from previous work (4). In order to test for this possibility, we performed a cluster analysis
220 on the predicted disease-associated epitopes bound by HLA-B (N = 132) and analyzed
221 cluster-specific motifs and HLA allele binding patterns. Intriguingly, among the ten most
222 dominant epitope clusters, each exhibiting a distinct peptide motif, nine were defined by
223 multiple HLA-B alleles (**Fig. 4**), some of them even belonging to different supertypes (*SI*

224 **Appendix, Table S7**). All of these clusters included both novel and previously described
225 epitopes, and three of them were defined by both risk- and protection-conferring alleles.
226 Furthermore, all HLA variants bound peptides of multiple dominant clusters; e.g. B*57:01
227 is associated with 3 dominant clusters, each showing a distinct peptide motif, but all
228 showing a strong preference for amino acid ‘W’ at anchor position 9 (**Fig. 4**). Overall, the
229 cluster analysis shows that our PepWAS approach identifies groups of peptides with
230 distinct motifs that are different from HLA variant-specific binding motifs (see also
231 schematic in **Fig. 1**). Generally, the 24 disease-associated epitopes predicted to be bound
232 by HLA-B*57:01 (but some of these also by other alleles), accounted for the highest level
233 of variation in spVL, even though they derived from 5 different HIV-1 proteins (**Fig. 3B,**
234 **C and SI Appendix, Fig. S8**). One of these epitopes, the well characterized HIV-1 *Gag*
235 epitope ISPRTLNAW (belonging to the dark purple cluster in **Fig. 4**), slightly exceeded
236 the effect of all other epitopes (**Fig. 3B**), in concordance with experimental evidence (17).
237 Other HLA-B alleles, including the B*08, B*44, and B*51 types, were also included in our
238 dataset, and their predicted epitope repertoires roughly followed their disease-association
239 known from previous studies (**SI Appendix, Fig. S1B**; allele-specific associations and
240 number of bound peptides are given in **SI Appendix, Table S4**).

241 Mechanistically, a negative association between predicted HIV-1 epitopes and viral load is
242 intuitive and likely resulting from the peptides’ immunodominant role in CTL response
243 and their escape mutations leading to significant fitness costs for the virus. However, a
244 number of the predicted HIV-1 epitopes exhibited a positive association with viral load,
245 indicating that they confer lower disease protection relative to the bulk of the peptides.
246 They likely represent peptide variants that fail to elicit an efficient CTL response or can
247 readily mutate with negligible fitness effects, thus allowing viral escape from HLA
248 presentation at no cost for the virus. Indeed, the most risk-associated predicted *Vpu* epitope,
249 IPIVAIVAL (**SI Appendix, Fig. S8F**; belonging to the largest, grey cluster in **Fig. 4**),
250 includes an anchor residue that exhibits significant variation in primary HIV-1 clones and
251 is involved in mediating immune-evasion through down-regulation of HLA-C (18), whose
252 high expression has been implicated in HIV control (19). The lack of significant
253 associations between predicted HLA-C bound epitopes and viral load in our analysis might
254 indicate that previously observed viral control associated with HLA-C is not mediated

255 through specific peptide presentation of HLA-C. However, more research is required to
256 fully understand the role of HLA-C in viral control (18).

257 So far, our analysis was based on the HIV-1 genome reference sequence. Though widely
258 used for research, focusing on this sequence accession may restrict our findings. We thus
259 repeated the entire analysis using the HIV-1 proteome consensus sequence from the Los
260 Alamos database, which incorporates major variation across different HIV-1 strains. The
261 results remained qualitatively the same (*SI Appendix, Fig. S9 and Table S8*). However,
262 HIV is well known to exhibit substantial within-host evolution (20, 21) and it is easily
263 conceivable that the ability of a patient's HLA variants to bind HIV epitopes is significantly
264 affected by genetic variation in the patient's HIV population (22). We therefore also
265 analyzed patient-specific autologous HIV-1 sequence information, which was available for
266 a small subset of patients, covering 8 of the 10 HIV-1 proteins (*SI Appendix, Table S6*).
267 For 4 of the 8 proteins (*Gag, Pol, Vif* and *Nef*) we found that the proportion of variation in
268 spVL associated with HLA-bound epitope repertoires changed when predicting epitopes
269 from autologous sequences instead of from the reference sequence. In all 4 cases, the
270 variation associated with predicted autologous epitopes was higher than when using their
271 homologs from the reference sequence (*SI Appendix, Table S6*), suggesting that our
272 PepWAS approach might be able to explain more variation in spVL than a standard GWAS
273 if autologous sequences were available for a larger fraction of infected individuals.

274 PepWAS relies on computational algorithms for the prediction of binding affinities
275 between HLA variants and peptides, and is thus inherently limited by their accuracy and
276 specificity. For instance, the empirical data used to train currently established HLA class I
277 algorithms contains mainly 9mer peptides, even though HLA class I molecules can
278 occasionally bind slightly shorter or longer peptides. Such peptides might therefore be
279 missed by current prediction algorithms. On the other hand, the current setup does in fact
280 identify 9mer cores of larger known epitopes. For instance, the here predicted protective
281 9mer *Gag* epitope 'STLQEIQGW' resides within the previously described 10mer *Gag*
282 epitope TW10 (**Fig. 3B**). Furthermore, this limitation is likely to be alleviated as more
283 training data is becoming available.

284 Overall, our findings reveal a functional basis of the robustly established association
285 between HLA genes and HIV-1 infection outcome. We show that both quantity and quality

286 of HLA-bound HIV epitopes contribute to controlling a patient's viral load. Our data also
287 suggests a more important role for *Env* protein-derived epitopes than previously thought.
288 Ultimately, our PepWAS approach of combining computational HLA-specific epitope
289 prediction with disease phenotype validation provides a promising avenue for
290 identification and prioritization of novel epitopes. As such, it complements existing
291 empirical essays for the development of targeted therapy. Noteworthy, by involving a
292 functional layer (peptide binding), the PepWAS approach enables the detection of disease-
293 relevant properties that are shared among several genetic variants (overlap in peptide
294 binding among HLA alleles). Such shared properties would be undetectable by GWAS,
295 because of its focus on distinct genetic variants instead of function, and should therefore
296 lead to higher sensitivity in the PepWAS approach compared to GWAS. Furthermore, the
297 PepWAS approach allows to account for individual variation in the pathogen proteome if
298 autologous sequence information is available, potentially further increasing sensitivity. As
299 such it may be applied to any HLA-associated complex disease.
300

301 **Material and Methods:**

302 For detailed information on Material and Methods see *SI Appendix* Supporting Methods
303 available online.

304

305 Samples and Genotype data:

306 We analyzed HLA genotype data and set point viral load (spVL) measurements of 6,311
307 subjects chronically infected with HIV-1. The original data and thorough quality check
308 are described in detail in McLaren *et al.* (4) and explained briefly in Supporting Methods.

309

310 HLA binding affinity for HIV-1 epitopes:

311 We used the NCBI accession NC_001802.1 as the reference sequence for the HIV-1
312 proteome (M group subtype B). The algorithm NetMHCcons-1.1 was used to predict HLA
313 allele-specific binding affinities for all 9mer peptides generated from the entire HIV-1
314 proteome, applying the default affinity rank threshold for ‘strongly bound’ peptides (rank
315 < 0.5).

316

317 Association with viral load:

318 The association of an allele or a peptide with viral load (spVL) was calculated using a
319 linear regression model corrected for population covariates following McLaren *et al.* (4).
320 Covariates included the first five principle components of SNP variation and the cohort
321 identity (all adopted from McLaren *et al.* (4)). Variation in viral load attributable to a
322 given variable (allele or epitope) was calculated as the difference between adjusted-R²
323 values of the model with variable and covariates and the model with covariates only,
324 following McLaren *et al.* (4). The variable’s regression coefficient was used as the
325 measure of its effect on viral load.

326

327 Clustering of HLA-B-specific predicted epitopes:

328 Position-associated entropy was calculated for all HLA-B-bound disease-associated
329 epitopes (N = 132) and used for visualization in a non-metric multidimensional scaling
330 plot as well as for density-based clustering.

331

332 HLA binding of peptides from autologous HIV-1 sequences:

333 We analyzed autologous HIV-1 sequences from Bartha et al. (23). Autologous sequences
334 were available for 8 of 10 HIV-1 proteins (only Gp41 segment for *Env*), and only for a
335 small subset of patients in our cohorts (*SI Appendix, Table S6*).

336

337 **Acknowledgements:**

338 Patient and HIV sequence data was collected and generously provided by the International
339 Collaboration for the Genomics of HIV. This project was supported by the Emmy Noether
340 Programme of the Deutsche Forschungsgemeinschaft (DFG grant LE 2593/3-1 to T.L.L.).
341 Furthermore, this project has been funded in whole or in part with federal funds from the
342 Frederick National Laboratory for Cancer Research, under Contract No.
343 HHSN261200800001E. The content of this publication does not necessarily reflect the
344 views or policies of the Department of Health and Human Services, nor does mention of
345 trade names, commercial products, or organizations imply endorsement by the U.S.
346 Government. This Research was also supported in part by the Intramural Research Program
347 of the NIH, Frederick National Lab, Center for Cancer Research.

348

349

350 **Competing interests:**

351 The authors declare no conflict of interest.

352

353 **References:**

354

- 355 1. Goulder PJR, Walker BD (2012) HIV and HLA Class I: An Evolving Relationship.
356 *Immunity* 37(3):426–440.
- 357 2. Fellay J, et al. (2007) A whole-genome association study of major determinants for
358 host control of HIV-1. *Science* (80-) 317(5840):944–947.
- 359 3. Pereyra F, et al. (2010) The major genetic determinants of HIV-1 control affect HLA
360 class I peptide presentation. *Science* (80-). doi:10.1126/science.1195271.
- 361 4. McLaren PJ, et al. (2015) Polymorphisms of large effect explain the majority of the
362 host genetic contribution to variation of HIV-1 virus load. *Proc Natl Acad Sci*
363 112(47):14658–14663.
- 364 5. Mellors JW, et al. (1996) Prognosis in HIV-1 infection predicted by the quantity of
365 virus in plasma. *Science* (80-) 272(5265):1167–1170.
- 366 6. Karosiene E, Lundegaard C, Lund O, Nielsen M (2012) NetMHCcons: A consensus
367 method for the major histocompatibility complex class I predictions.
368 *Immunogenetics* 64(3):177–186.
- 369 7. Zhang H, Lundegaard C, Nielsen M (2009) Pan-specific MHC class I predictors: A
370 benchmark of HLA class I pan-specific prediction methods. *Bioinformatics*
371 25(1):83–89.
- 372 8. Rooney MS, Shukla SA, Wu CJ, Getz G, Hacohen N (2015) Molecular and genetic
373 properties of tumors associated with local immune cytolytic activity. *Cell* 160(1–
374 2):48–61.
- 375 9. Strønen E, et al. (2016) Targeting of cancer neoantigens with donor-derived T cell
376 receptor repertoires. *Science* (80-) 352(6291):1337–1341.
- 377 10. Košmrlj A, et al. (2010) Effects of thymic selection of the T-cell repertoire on HLA
378 class I-associated control of HIV infection. *Nature* 465(7296):350–354.
- 379 11. Martoglio B (1997) Signal peptide fragments of preprolactin and HIV-1 p-gp160
380 interact with calmodulin. *EMBO J* 16(22):6636–6645.
- 381 12. Yusim K, et al. (2009) HIV molecular immunology. *Los Alamos, New Mex Los*
382 *Alamos Natl Lab Theor Biol Biophys*:3–24.
- 383 13. Pereyra F, et al. (2014) HIV Control Is Mediated in Part by CD8+ T-Cell Targeting

384 of Specific Epitopes. *J Virol* 88(22):12937–12948.

385 14. Kiepiela P, et al. (2004) Dominant influence of HLA-B in mediating the potential
386 co-evolution of HIV and HLA. *Nature* 432(7018):769–775.

387 15. Borrow P, Lewicki H, Hahn BH, Shaw GM, Oldstone MBA (1994) Virus-Specific
388 CD8+ Cytotoxic T-Lymphocyte Activity Associated with Control of Viremia in
389 Primary Human Immunodeficiency. *68(9):6103–6110*.

390 16. Falk K, Rötzschke O, Stevanović S, Jung G, Rammensee HG (1991) Allele-specific
391 motifs revealed by sequencing of self-peptides eluted from MHC molecules. *Nature*
392 351(6324):290–296.

393 17. Llano A, Williams A, Olvera A, Silva-Arrieta S, Brander C (2013) Best-
394 Characterized HIV-1 CTL Epitopes: The 2013 Update. *HIV Mol Immunol* 2013:3–
395 25.

396 18. Apps R, et al. (2016) HIV-1 Vpu Mediates HLA-C Downregulation. *Cell Host*
397 *Microbe* 19(5):686–695.

398 19. Thomas R, et al. (2009) HLA-C cell surface expression and control of HIV/AIDS
399 correlate with a variant upstream of HLA-C. *Nat Genet* 41(12):1290–1294.

400 20. Cotton LA, et al. (2014) Genotypic and Functional Impact of HIV-1 Adaptation to
401 Its Host Population during the North American Epidemic. *PLoS Genet* 10(4).
402 doi:10.1371/journal.pgen.1004295.

403 21. Li G, et al. (2015) An integrated map of HIV genome-wide variation from a
404 population perspective. *Retrovirology* 12(1). doi:10.1186/s12977-015-0148-6.

405 22. Kawashima Y, et al. (2009) Adaptation of HIV-1 to human leukocyte antigen class
406 I. *Nature* 458(7238):641–645.

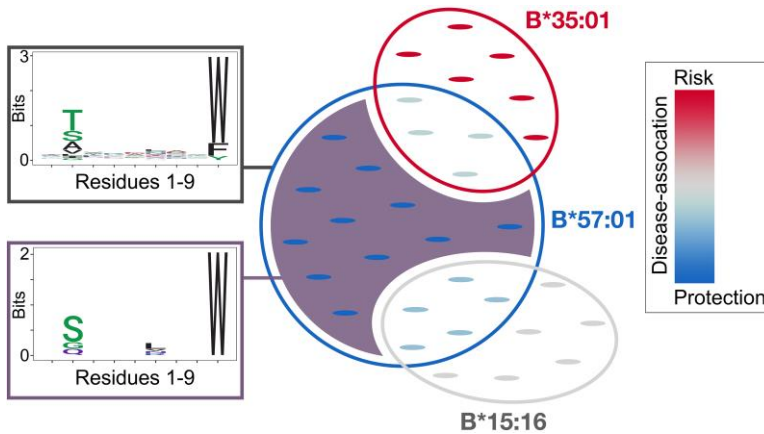
407 23. Bartha I, et al. (2013) A genome-to-genome analysis of associations between human
408 genetic variation, HIV-1 sequence diversity, and viral control. *Elife* 2013(2):1–16.

409 24. McLaren PJ, Carrington M (2015) The impact of host genetic variation on infection
410 with HIV-1. *Nat Immunol* 16(6):577–583.

411

412

413 **Figure legends:**

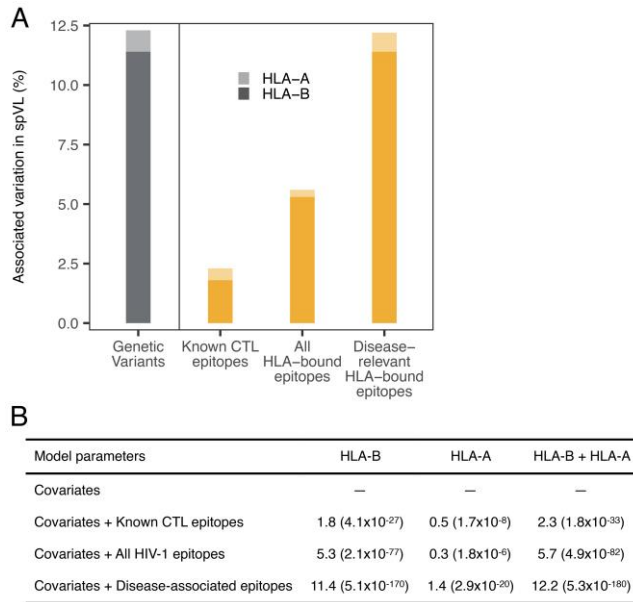


414

415 **Fig. 1. Schematic for determining peptide-specific associations through PepWAS.**

416 Disease-associated peptides are identified by integrating the different disease-associations
417 of the different HLA alleles that are predicted to bind them. Some peptides will only be
418 bound by one HLA allele, thus drawing their disease-association directly from the disease-
419 association of that allele (e.g. peptides in the purple shaded area, bound only by HLA-
420 B*57:01). However, many peptides will be bound by several HLA alleles, which can have
421 quite distinct, possibly even opposing disease associations (e.g. peptides in overlap of
422 *B57:01 and *B35:01). In this case, the disease-association of the peptide derives from the
423 disease-associations of each of the binding HLA alleles as well as their frequencies in the
424 dataset. The novel peptidome-wide association study (PepWAS) approach differentiates
425 these distinct sets of peptides and identifies both specific peptides and epitope motifs with
426 distinct disease-association (e.g. distinct motif of purple shaded peptides, corresponding to
427 the dark purple cluster in Fig. 5). Circles depict repertoires of peptides (small pointed ovals)
428 predicted to be bound by the given HLA allele. Overlap of circles defines sets of peptides
429 bound by both HLA alleles. Color of circles and peptides depicts disease-association of
430 corresponding HLA alleles and peptides, respectively, from blue (protective) to red (risk).
431 The number of peptides in this schematic does not correspond to the actual number of
432 peptides observed for these HLA alleles. In reality, the overlap among HLA alleles is
433 substantially more complex than depicted in this simplified schematic.

434

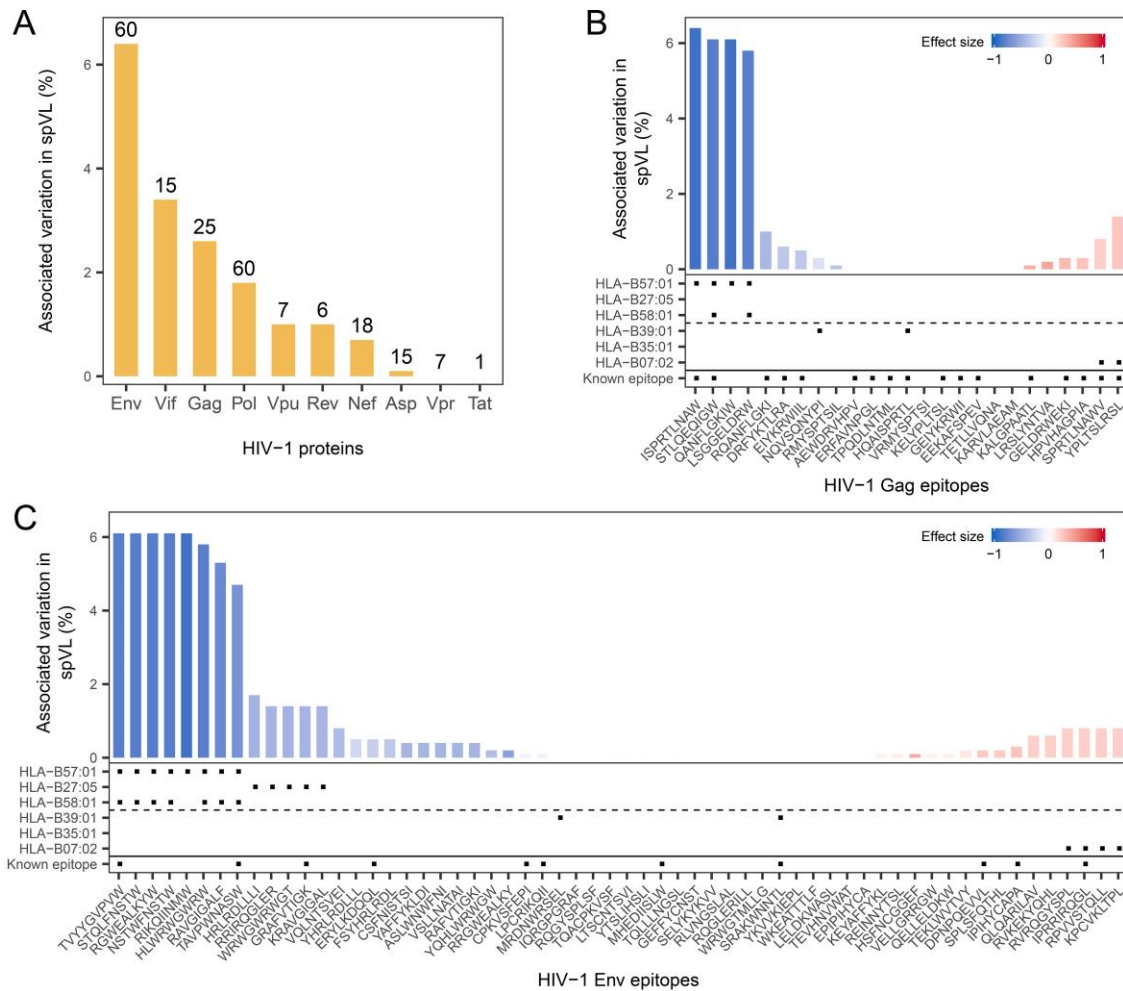


435

436 **Fig. 2. Variation in viral load associated with predicted epitope repertoires bound by**
 437 **HLA-B and HLA-A.**

438 Among HIV patients (N = 6,311), the proportion of variation (estimated as adjusted ΔR^2)
 439 in set point viral load (spVL) associated with the patient-specific number of predicted
 440 HLA-bound HIV-1 epitopes is shown separately for HLA-B and HLA-A, and for different
 441 epitope sets. **(A)** Previously, 11.4% and 0.9% of the variation in spVL had been associated
 442 with independent genetic variants in HLA-B and HLA-A, respectively (grey bars; data
 443 from ref. 4). Here we instead calculated the variation in spVL associated with individual
 444 HLA-bound HIV epitope repertoires (yellow bars), based on known CTL epitopes from
 445 Los Alamos HIV Molecular Immunology Database, all HLA-bound HIV epitopes, and
 446 only the disease-associated HIV epitopes (the latter corresponding to 99.2% of the variation
 447 previously associated with HLA genetic variation). **(B)** Variation associated with different
 448 sets of predicted epitopes. *P*-values (in parentheses) indicate the improvement over null
 449 model (covariates only: first five PCs and cohort group). Number of disease-associated
 450 predicted epitopes is 132 for HLA-B, and 74 for HLA-A, respectively.

451

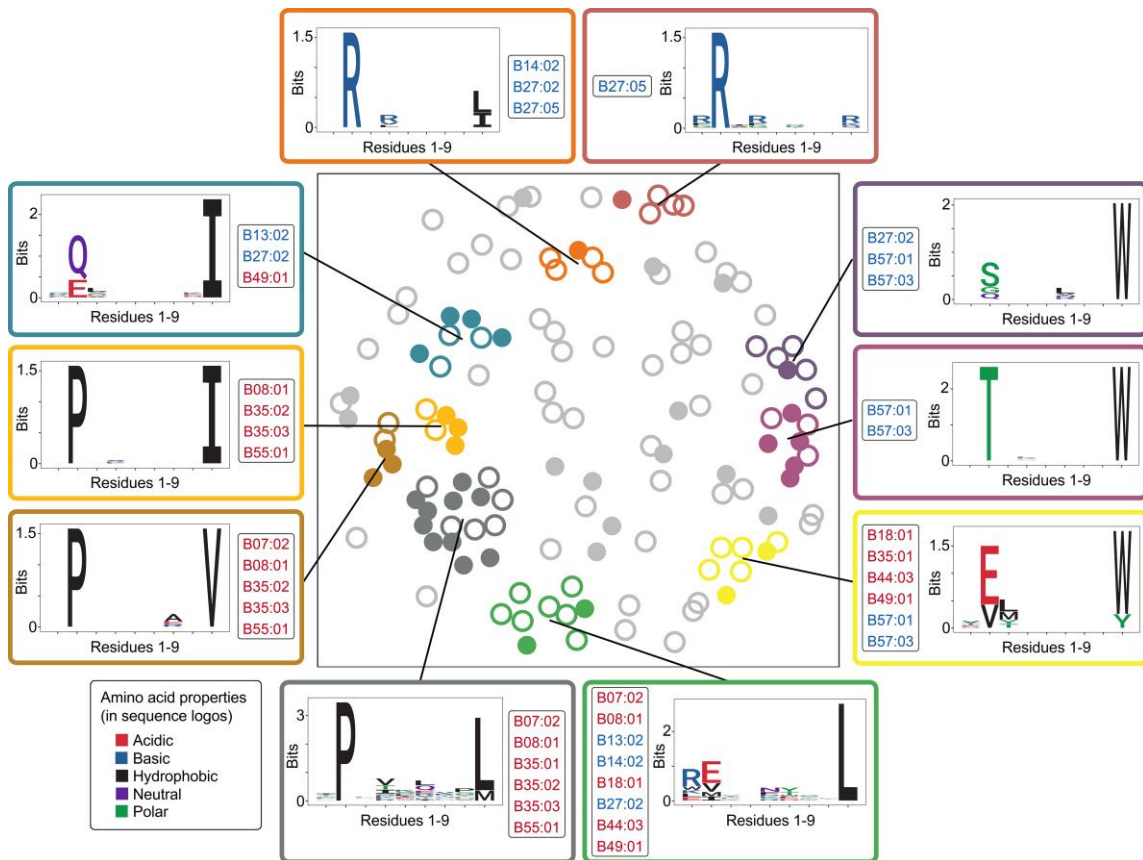


452

453 **Fig. 3. Epitope- and protein-specific association with viral load.**

454 (A) Percent of variation in spVL associated with all predicted epitopes of a given HIV-1
 455 protein. Absolute number of predicted HLA-B bound epitopes per protein is shown above
 456 the bars. (B-C) Predicted HLA-B-bound epitopes accounted for varied levels of variation
 457 in set point viral load (spVL). Height of the bar represents the fraction of variation in spVL
 458 associated with each epitope, while the color reflects each epitope's effect on spVL,
 459 ranging from protection (blue) to risk (red). Note that epitope effects are estimated
 460 separately and are thus not independent. *Gag* (B) and *Env* (C) proteins are shown as
 461 representative examples, together with information on predicted binding for 3 protective
 462 and 3 risk HLA-B alleles highlighted in a recent review (24) and whether peptides are
 463 known epitopes in Los Alamos HIV database. All other HIV-1 proteins are shown in *SI*
 464 *Appendix, Fig. S8.*

465



466

467 **Fig. 4. Clusters of disease-associated epitopes.**

468 Non-metric multidimensional scaling (NMDS) was used to visualize the pairwise distance
 469 between predicted HLA-B-bound disease-associated epitopes, which revealed 10 dominant
 470 clusters. Each circle represents an HLA-B bound disease-associated epitope (N = 132).
 471 Filled circles represent known CTL epitopes from the Los Alamos HIV Molecular
 472 Immunology Database (N = 45), while open circles represent previously uncharacterized
 473 disease-associated predicted epitopes. Cluster-specific motif and HIV-1 associated HLA-
 474 B alleles (N = 16) binding the cluster's epitopes are shown. The coloring of the allele names
 475 indicates disease-association of the specific alleles.