

Effects of thymic selection of the T-cell repertoire on HLA class I-associated control of HIV infection

Andrej Košmrlj^{1,2*}, Elizabeth L. Read^{1,3,4*}, Ying Qi⁵, Todd M. Allen¹, Marcus Altfeld¹, Steven G. Deeks⁶, Florencia Pereyra¹, Mary Carrington^{1,5}, Bruce D. Walker^{1,7} & Arup K. Chakraborty^{1,3,4,8}

Without therapy, most people infected with human immunodeficiency virus (HIV) ultimately progress to AIDS. Rare individuals ('elite controllers') maintain very low levels of HIV RNA without therapy, thereby making disease progression and transmission unlikely. Certain HLA class I alleles are markedly enriched in elite controllers, with the highest association observed for *HLA-B57* (ref. 1). Because HLA molecules present viral peptides that activate CD8⁺ T cells, an immune-mediated mechanism is probably responsible for superior control of HIV. Here we describe how the peptide-binding characteristics of HLA-B57 molecules affect thymic development such that, compared to other HLA-restricted T cells, a larger fraction of the naive repertoire of B57-restricted clones recognizes a viral epitope, and these T cells are more cross-reactive to mutants of targeted epitopes. Our calculations predict that such a T-cell repertoire imposes strong immune pressure on immunodominant HIV epitopes and emergent mutants, thereby promoting efficient control of the virus. Supporting these predictions, in a large cohort of HLA-typed individuals, our experiments show that the relative ability of HLA-B alleles to control HIV correlates with their peptide-binding characteristics that affect thymic development. Our results provide a conceptual framework that unifies diverse empirical observations, and have implications for vaccination strategies.

HIV infection leads to acute high level viraemia, which is subsequently reduced to a set-point viral load. Without therapy, most patients experience a subsequent increase in viral load, and ultimately the development of AIDS. Viraemia levels and time to disease vary widely, and the differences correlate with the expression of different HLA class I molecules (reviewed in ref. 2). Effector CD8⁺ T cells (CTLs) are implicated in viral control because T-cell antigen receptors (TCRs) on CD8⁺ T cells recognize complexes of viral peptides and class I HLA molecules presented on the surface of infected cells, and depletion of CD8⁺ T cells leads to increased viraemia in animal models of HIV infection³. We describe a feature of the HLA-B57-restricted CD8⁺ T-cell repertoire that contributes to enhanced control of viraemia.

Algorithms⁴ based on experimental data predict whether a particular peptide will bind to a given HLA molecule⁵. We tested four predictive algorithms against available experimental data on peptide binding to diverse HLA molecules and found that, in most cases, they are highly accurate (Supplementary Fig. 1 and Supplementary Table 1). For example, predictions using the best algorithm for HLA-B*5701 were 97% accurate. Using these algorithms, we computed the fraction of peptides derived from the human proteome⁶ that bind to various HLA

molecules. Of the $\sim 10^7$ unique peptide sequences, only 70,000 are predicted to bind to HLA-B*5701, and 180,000 bind to HLA-B*0701 (an allele that is not protective against HIV). Essentially identical results were obtained for randomly generated peptides (data not shown). The protective allele in macaques, Mamu-B*17, also binds fewer self peptides than other Mamu molecules for which data are available (Mamu-B*17 binds 4, 6 and 13 times fewer self peptides than Mamu-A*11, Mamu-A*01 and Mamu-A*02, respectively; Supplementary Table 1).

The intrinsic differences in self-peptide binding among HLA molecules are important during T-cell repertoire development. Immature T cells are exposed to diverse host-derived peptide–HLA complexes presented in the thymus. As fewer self peptides are able to bind to HLA-B*5701 (and Mamu-B*17) molecules, a smaller diversity of self-peptide TCR contact sequences will be encountered by HLA-B*5701/Mamu-B*17-restricted T cells in the thymus (Supplementary Discussion 1).

The diversity of self peptides presented in the thymus shapes the characteristics of the mature T-cell repertoire. Experiments^{7,8} and theoretical studies^{9,10} show that T cells that develop in mice with only one type of peptide in the thymus are more cross-reactive to point mutants of peptide epitopes that they recognize than T cells from mice that express diverse self peptides. Thus, by encountering fewer self peptides during thymic development, HLA-B57-restricted CD8⁺ T cells should be more cross-reactive to point mutants of targeted viral peptides.

We carried out *in silico* thymic selection experiments to test this hypothesis. We chose an HLA-dependent number of thymic self peptides, each with amino acids of the TCR contact residues picked according to the frequency with which they appear in the human proteome^{6,9}. A diverse set of immature CD8⁺ T cells (thymocytes) was generated by choosing the sequences of their peptide contact residues in the same way, and by varying the TCR–HLA interactions. A thymocyte emerges from the thymus as a mature CD8⁺ T cell if its TCR binds to at least one self-peptide–major histocompatibility complex (pMHC; human MHC is called HLA) molecule with an affinity that exceeds the positive selection threshold, and does not interact with any pMHC more strongly than the negative selection threshold. Using a computational model^{9,10} in the class of 'string models'¹¹, we assessed the affinity of TCR–self-peptide–HLA complexes (Methods) to determine which T cells survive positive and negative selection, and become a part of the mature repertoire. Our qualitative results are independent of the parameters used to determine these interaction strengths (Supplementary Figs 2 and 3)^{9,10}.

¹Ragon Institute of MGH, MIT and Harvard, Boston, Massachusetts 02114, USA. ²Department of Physics, Massachusetts Institute of Technology, Cambridge, Massachusetts 02139, USA. ³Department of Chemical Engineering, Massachusetts Institute of Technology, Cambridge, Massachusetts 02139, USA. ⁴Department of Chemistry, Massachusetts Institute of Technology, Cambridge, Massachusetts 02139, USA. ⁵Cancer and Inflammation Program, Laboratory of Experimental Immunology, SAIC-Frederick, Inc., NCI-Frederick, Frederick, Maryland 21702, USA. ⁶University of California, San Francisco, California 94110, USA. ⁷Howard Hughes Medical Institute, Chevy Chase, Maryland 20815, USA. ⁸Department of Biological Engineering, Massachusetts Institute of Technology, Cambridge, Massachusetts 02139, USA.

*These authors contributed equally to this work.

The mature T cells that emerged from these *in silico* thymic selection experiments were then computationally challenged by a viral peptide (that is, not seen in the thymus) bound to the same HLA type. T cells that recognize this peptide–HLA complex were obtained by assessing whether the interaction strength exceeded the negative selection threshold (shown to be equal to the recognition threshold in mouse models¹²); qualitative results are invariant if the recognition threshold is not much weaker than that corresponding to negative selection (Supplementary Fig. 3). Cross-reactivity of these T cells was then determined *in silico* by mutating each TCR contact residue of the peptide to the other 19 possibilities. Sites on the viral peptide were called ‘important contacts’ if half the mutations therein abrogated recognition by T cells that target this epitope. The frequency of the number of important contacts in viral peptides that determine T-cell recognition was obtained by repeating this procedure 1,000 times with different choices of thymocytes and self and foreign peptides.

Our calculations predict that a T-cell repertoire restricted by an HLA molecule such as HLA-B*5701, which presents fewer self peptides in the thymus, has a higher frequency of occurrence of T cells that recognize viral peptides through smaller numbers of important contacts (Fig. 1a). In contrast, the frequency of occurrence of T cells

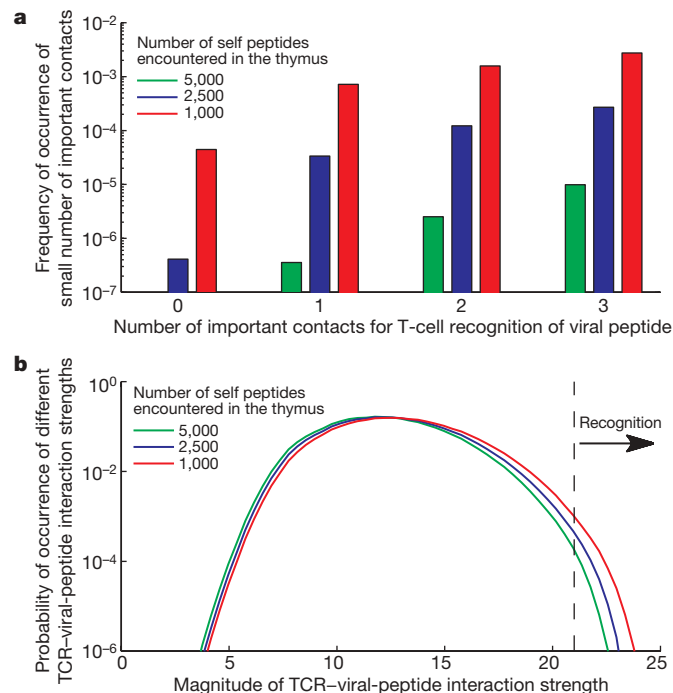


Figure 1 | Thymic selection against fewer self peptides leads to a more cross-reactive T-cell repertoire. **a**, Histogram of the frequency with which T cells recognize viral peptides (that is, not seen in the thymus) through only a small number (0, 1, 2, 3) of important contacts is shown for three T-cell repertoires that developed with different numbers of self-peptide–HLA complexes in the thymus. Important contacts were determined by making single point mutations. If the TCR–peptide–HLA interaction is sufficiently strong, no single point mutation can abrogate recognition, resulting in zero important contacts. A higher frequency of occurrence of a small number of important contacts indicates a more cross-reactive T-cell repertoire because only mutations at these contacts are likely to abrogate recognition. The frequency with which T cells recognize viral peptides through many significant contacts (greater than four) is larger for T-cell repertoires restricted by HLA alleles that present more self peptides in the thymus (not shown). **b**, The probability that a TCR binds to viral peptides with a certain interaction strength is shown for three T-cell repertoires (as in **a**). A particular TCR recognizes a viral peptide when the binding strength exceeds the recognition threshold (dotted black line). Members of a T-cell repertoire selected against fewer self peptides are more likely to recognize a viral peptide. The model we used describes qualitative trends robustly^{9,10} (Methods), but is not meant to be quantitatively accurate.

that recognize viral peptides through many important contacts is larger for repertoires restricted by HLA alleles that present a greater diversity of self peptides in the thymus (data not shown for >four contacts). Mutations at sites different from the important contacts do not affect binding strength substantially. Therefore, when the interaction between peptide–HLA and TCR is mediated by fewer important contacts, a larger number of possible point mutations of the peptide do not affect peptide recognition, thereby making the T cells more cross-reactive to mutants that arise. Thus, the HLA-B57-restricted T-cell repertoire is expected to be more cross-reactive to mutants of targeted viral peptides than repertoires restricted by HLA alleles that present a greater diversity of self peptides.

Our computational models give this qualitative mechanistic insight, but do not provide quantitative estimates of the extent of this enhanced cross-reactivity of T cells. However, compelling experimental data¹³ has shown that the effect revealed by our studies is important in humans. Peripheral blood mononuclear cells from patients expressing HLA-B57 contained CTLs that were more cross-reactive to various HIV epitopes and their point mutants than those of HLA-B8-positive patients. HLA-B8 is associated with rapid progression to disease¹³, and the most accurate algorithm for peptide binding suggests that the HLA-B8 molecule binds a greater diversity of self peptides than HLA-B57 (Supplementary Fig. 4 and Supplementary Table 1). Other experimental studies also show that patients expressing HLA-B57 cross-recognize point mutants of the dominant epitope and use more public TCRs^{14,15}.

Next, we computed interaction strengths between diverse viral peptides and members of T-cell repertoires restricted by HLA molecules that present differing numbers of self peptides in the thymus. This allowed us to obtain the probability with which a randomly picked T-cell clone and viral peptide will interact sufficiently strongly for recognition to occur. The results (Fig. 1b) indicate that a typical CD8⁺ T cell restricted by an HLA molecule such as HLA-B*5701, which presents fewer peptides in the thymus, has a higher probability of recognizing a viral epitope compared to a T cell restricted by other HLA molecules. Thus, more HLA-B*5701-restricted T cell clones are likely to recognize a viral epitope, making effective precursor frequencies higher in an HLA-B*5701-restricted repertoire (a strong predictor of response magnitude¹⁶). A greater precursor frequency for viral epitopes in the naive repertoire restricted by HLA-B57 is indicated by experimental results showing that HLA-B*5701 contributes the most to acute-phase CTL responses of all HLA alleles tested¹⁷.

The results in Fig. 1 stem from the constraint that thymocytes must avoid being negatively selected by each self-peptide–HLA complex encountered during development in the thymus. T cells expressing TCRs with peptide contact residues composed of amino acids that interact strongly with other amino acids (for example, charged residues, flexible side chains) have a high probability of binding to a self peptide strongly. The greater the diversity of self peptides presented in the thymus, the higher the chance that a TCR with such peptide contact residues will encounter a self peptide with which strong interactions will result in negative selection. Thus, as the diversity of self peptides presented in the thymus increases, the peptide contact residues of TCRs in the mature T-cell repertoire are increasingly enriched in weakly interacting amino acids (Supplementary Fig. 5). T cells bearing TCRs with weakly interacting peptide contact residues recognize viral peptides by means of several moderate interactions, making many contacts important for recognition. In contrast, TCRs with peptide contact residues containing strongly interacting amino acids are more likely to recognize viral peptides through a few important contacts mediated by these residues, making recognition cross-reactive to mutations at other peptide sites. These mechanistic insights are supported by experimental results^{7,9} (Supplementary Discussion 2).

By studying a model of host–pathogen dynamics that builds on past models of host–HIV interactions^{18–20}, we explored the consequences

of the HLA-B57-restricted CD8⁺ T-cell repertoire having a higher precursor frequency for viral peptides and being more cross-reactive to point mutants of targeted epitopes on the control of HIV. Because of the importance of immune control exerted by CD8⁺ T cells^{17,21}, we focused on the interaction between a mutating virus quasispecies and epitope-directed, variably cross-reactive, host CTL responses.

The essential features of the model are depicted in Fig. 2a (details in Methods). The virus is modelled as a number of epitopes consisting of strings of amino acids, and new viral strains (point mutations of epitopes), which differ in replicative fitness, arise over the course of infection. For each individual, an HLA-dependent CD8⁺ T-cell repertoire was chosen. To mimic the results obtained from our thymic selection calculations (Fig. 1b), more or less cross-reactive repertoires were chosen (Supplementary Fig. 6) to represent HLA-B57-restricted T cells and those restricted by other HLAs, respectively. Infection rates were limited by target CD4⁺ T cells, and CTL contraction and memory were included. Other dynamic models were studied, including one that does not incorporate target cell limitation or CTL contraction. Our qualitative results about the effects of cross-reactivity are robust to variations in parameters and model assumptions (Supplementary Figs 7–16).

We find that individuals with a more cross-reactive CTL repertoire control viral loads better during the acute phase of the infection (Fig. 2b). This is in agreement with findings in simian immunodeficiency virus (SIV)-infected rhesus macaques²², where the number of cross-reactive TCR clones negatively correlates with viral load. Our simulations show that a larger number of CTL clones in a more cross-reactive T-cell repertoire recognize epitopes from the infecting viral strain (Fig. 2c). This is because the predicted higher precursor frequency for viral epitopes (Fig. 1b) leads to a greater response magnitude (as in mouse models¹⁶). This conclusion is supported by data showing that in people with a protective HLA allele, the initial T-cell response to HIV is dominated by T cells restricted by the protective HLA and not those restricted by other HLAs expressed¹⁷. Our simulations also show that enhanced cross-reactivity of the T-cell repertoire leads to greater immune pressure on the emergent viral mutants by individuals expressing HLA-B57 compared to those with T cells restricted by HLA molecules that bind more types of self peptides. The stronger immune pressure on infecting and emerging viral strains results in superior control of viral load. Thus, we predict that HIV-infected individuals with HLA alleles that bind fewer self peptides are more likely to control viral loads to low values.

To test this prediction, we studied two large HLA-typed cohorts: 1,110 controllers with less than 2,000 HIV particles ml⁻¹ and 628 progressors (or non-controllers) with viral loads exceeding 10⁴ ml⁻¹ (Methods). From these data, we obtained the odds ratio (OR) for individual HLA alleles. People with HLA alleles associated with an OR value greater or less than one are more likely to be progressors or controllers, respectively. We focused on HLA-B alleles because they are associated with control of HIV²³. Of 40 HLA-B alleles that were studied, significant results (*P* value < 0.05) were obtained for five HLA-B alleles (Supplementary Table 2) and peptide-binding data are available for four of them. In support of our predictions, those HLA-B alleles associated with higher OR values also bind more self peptides (Fig. 3).

Superior control of viral load due to the greater precursor frequency and cross-reactivity of those T-cell repertoires restricted by HLA molecules that bind to few self peptides (for example, HLA-B57)

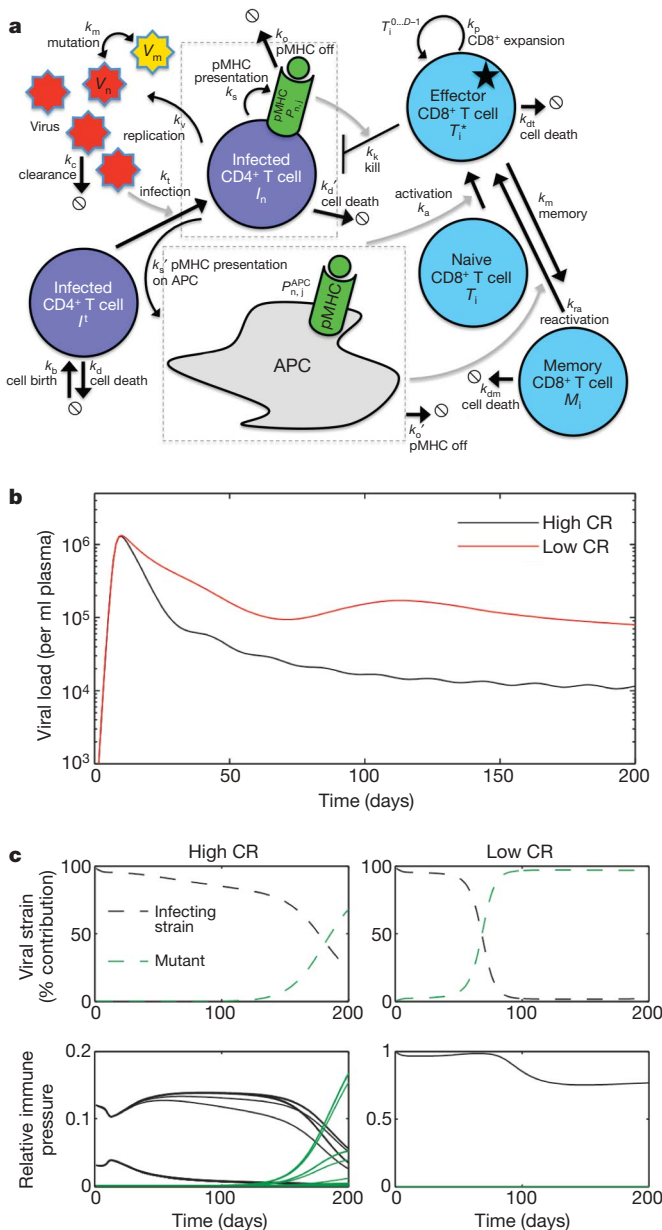


Figure 2 | Model of host–pathogen interactions shows superior viral control by cross-reactive CD8⁺ T-cell repertoires. **a**, Dynamic model: the virus mutates, infects limited-target CD4⁺ T cells, and is cleared. Infected CD4⁺ T cells produce more free virus and die. Infected cells present viral peptides in complex with HLA molecules (until peptides unbind from HLA). Activated CD8⁺ T cells produced by recognition of viral epitopes on antigen-presenting cells (APCs) proliferate and differentiate into effector CTLs. CTLs kill infected cells bearing cognate peptide–HLA complexes, and turn into memory cells that are activated after re-exposure to antigen. **b**, Simulated HIV viral loads versus time for different cross-reactivities (CR) of the CD8⁺ T-cell repertoire. Black curve, high cross-reactivity; red curve, low cross-reactivity. Each curve is averaged over 500 simulations (each simulation represents a person). The model shows a reduced set-point viral load for people with a more cross-reactive T-cell repertoire. Other models of host–pathogen dynamics show similar effects of T-cell cross-reactivity (Supplementary Figs 7 and 8). **c**, Virus diversity and immune pressure for representative people (that is, representative simulations) with high cross-reactivity (left) and low cross-reactivity (right) of CD8⁺ T-cell repertoires. Top panels show the relative population sizes of two dominant viral strains: the infecting strain (black), and an emerging, less fit strain (green) (other less populous viral strains are not shown). For people with a more cross-reactive T-cell repertoire, the emergent mutant strain only begins to dominate the infecting strain after 175 days, whereas for low cross-reactivity the mutant increases to nearly 100% of the viral population within 100 days after infection. Bottom panels show the relative immune pressure, defined as the rate of killing of an infected cell (see equation (4), Methods), imposed on each viral strain by different CD8⁺ T-cell clones. Each curve represents the relative immune pressure exerted on that viral strain by a particular T-cell clone. For people with a more cross-reactive T-cell repertoire, several T-cell clones exert immune pressure on both the infecting and emergent strains. For people with a low-cross-reactivity T-cell repertoire, the emergent strain is not recognized, and thus escapes.

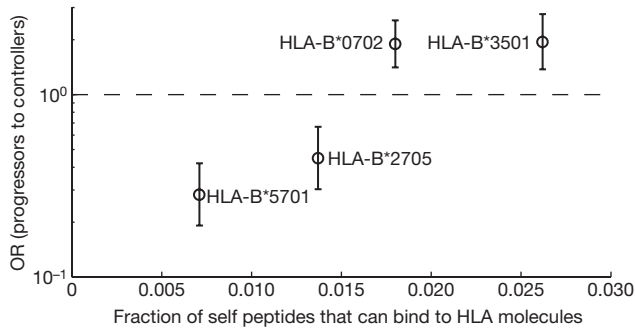


Figure 3 | HLA-B alleles associated with greater ability to control HIV correlate with smaller self-peptide binding propensities. The odds ratio (OR) for an allele is defined as: $\frac{p_w/p_{wo}}{c_w/c_{wo}}$, where p_w and p_{wo} are the numbers of individuals in the progressor cohort with and without this HLA, respectively; and c_w and c_{wo} are the numbers of individuals in the controller cohort with and without this HLA, respectively. This definition suggests that the OR measures the likelihood of an allele being correlated with progressors versus controllers, with an OR greater than one indicating association with the progressor cohort. The fraction of peptides derived from the human proteome that bind to a given HLA allele was determined using the most accurate predictive algorithms (Methods and Supplementary Table 1). Compared to experimental data, the predictive algorithms for peptide binding by HLA-B*3501 are less accurate than algorithms for the other three alleles (Supplementary Fig. 17 and Supplementary Table 1); the number reported here for HLA-B*3501 using the most accurate algorithm underestimates the binding fraction. The error bars represent the 95% confidence intervals for OR. The dotted line corresponds to equal odds for an allele being associated with progressors and controllers.

should also confer protection against diseases caused by other fast-mutating viruses. Indeed, HLA-B57 is protective against hepatitis C virus (HCV)²⁴, another highly mutable viral disease in which CD8⁺ T cells are important. Also, HLA-B8, which binds a greater diversity of self peptides, is associated with faster disease progression in HCV²⁵ and HIV¹³. Thus, the correlation between the diversity of peptides presented in the thymus during T-cell development and control or progression of disease may be general.

Undoubtedly, many complex factors influence the relationship between HLA type and disease outcome. The effect of the new factor we have identified should be greatest for HLA molecules that bind relatively few (for example, HLA-B57) or many (for example, HLA-B7, -B35, -B8) self peptides. The strong association of HLA-B27—which binds an intermediate number of self peptides (twice as many as HLA-B57)—with viral control indicates that, in this case, the effects of T cell cross-reactivity are reinforced by this molecule binding HIV epitopes that are subject to very strong structural constraints.

Our results also point to a mechanistic explanation for as yet unexplained associations between HLA alleles that confer protection against HIV and autoimmune diseases. T cells restricted by HLA alleles that bind to few self peptides are subject to less stringent negative selection in the thymus, and should therefore be more prone to recognizing self peptides. Indeed, HLA-B57 has been associated with autoimmune psoriasis²⁶ and hypersensitivity reactions²⁷. Enhanced cross-reactivity of HLA-B27-restricted T cells and other unique properties of this molecule (misfolding, homodimers²⁸) probably contribute to the enhanced risk of autoimmunity associated with this allele²⁹.

Our results shed light on another intriguing observation; acutely infected patients with low viral loads (and protective HLAs) tend to target an immunodominant epitope that makes a larger relative contribution to the total CTL response as compared to individuals presenting with higher levels of viraemia³⁰. This is counterintuitive as the most protective responses appear most focused, rather than broadly distributed over many epitopes. We calculated how viral load correlates with the number of CTLs responding to the immunodominant epitope divided by the total number of CTLs activated by the virus (a

quantity analogous to relative contribution³⁰). Mirroring experimental data, HLA alleles that restrict a more cross-reactive repertoire and are more protective also make a larger relative contribution (Supplementary Fig. 13). This result unifies the idea of both a broad and a focused response. The more cross-reactive repertoire targets more epitopes and emergent mutants, but a larger number of clones also recognize the dominant epitope (Fig. 2c).

Cross-reactive T cells are rare in people with HLA alleles that present more self peptides in the thymus than the B57 allele, but they do exist. Our results suggest that a T-cell vaccine for a diverse population must aim to activate these rare cross-reactive T cells that also target epitopes from a conserved region of the HIV genome (like HLA-B57 Gag epitopes). This will enable robust responses to infecting and mutant strains until a strain with low replicative fitness emerges, enhancing control of viral load.

METHODS SUMMARY

Predictive algorithm tools for peptide binding to HLA and Mamu molecules were obtained from the Immune Epitope Database (IEDB)⁴ and were used to predict the fraction of bound peptide derived from the human and macaque proteomes⁶. Accuracies of these tools were tested on experimental data obtained from the IEDB⁴. To assess the effects of thymic selection on TCRs restricted by different MHC molecules (HLA or Mamu), we used a computational model of thymic selection described in Methods (and previously^{9,10}).

To explore host–pathogen dynamics, we constructed a small model of the HIV virus with distinct epitopes and sequence diversity, based in part on past work^{18–20}. We carried out numerical simulations of ordinary differential equation models, shown schematically in Fig. 2a and Supplementary Fig. 7. Parameters and their justification are given in Supplementary Tables 3 and 4 and in the Supplementary Methods. To explore cross-reactivity, we varied the distribution of pairwise-interaction free energies of TCR–pMHC contacts. Our goal was not to obtain precise numbers, but to examine the qualitative effects of variation in repertoire cross-reactivity on virus control. Qualitative results are robust to variations in parameters and assumptions (Supplementary Figs 8–16).

HLA-typed cohorts of people of diverse races were divided into HIV controllers and HIV non-controllers, and analysed for HLA association with the ability to control HIV. The results (Fig. 3 and Supplementary Table 2) were adjusted for the effects of HLA-B*0702, HLA-B*3501, HLA-B*2705 and HLA-B*5701.

Full Methods and any associated references are available in the online version of the paper at www.nature.com/nature.

Received 13 October 2009; accepted 11 March 2010.

Published online 5 May 2010.

- Migueles, S. A. *et al.* HLA B*5701 is highly associated with restriction of virus replication in a subgroup of HIV-infected long term nonprogressors. *Proc. Natl Acad. Sci. USA* **97**, 2709–2714 (2000).
- Deeks, S. G. & Walker, B. D. Human immunodeficiency virus controllers: mechanisms of durable virus control in the absence of antiretroviral therapy. *Immunity* **27**, 406–416 (2007).
- Jin, X. *et al.* Dramatic rise in plasma viremia after CD8⁺ T cell depletion in simian immunodeficiency virus-infected macaques. *J. Exp. Med.* **189**, 991–998 (1999).
- Peters, B. *et al.* The immune epitope database and analysis resource: from vision to blueprint. *PLoS Biol.* **3**, e91 (2005).
- Rao, X., Fontaine Costa, A. I. C. A., van Baarle, D. & Kesmir, C. A comparative study of HLA binding affinity and ligand diversity: implications for generating immunodominant CD8⁺ T cell responses. *J. Immunol.* **182**, 1526–1532 (2009).
- Hubbard, T. J. *et al.* Ensembl 2009. *Nucleic Acids Res.* **37**, D690–D697 (2009).
- Huseby, E. S., Crawford, F., White, J., Marrack, P. & Kappler, J. W. Interface-disrupting amino acids establish specificity between T cell receptors and complexes of major histocompatibility complex and peptide. *Nature Immunol.* **7**, 1191–1199 (2006).
- Huseby, E. S. *et al.* How the T cell repertoire becomes peptide and MHC specific. *Cell* **122**, 247–260 (2005).
- Košmrlj, A., Jha, A. K., Huseby, E. S., Kardar, M. & Chakraborty, A. K. How the thymus designs antigen-specific and self-tolerant T cell receptor sequences. *Proc. Natl Acad. Sci. USA* **105**, 16671–16676 (2008).
- Košmrlj, A., Chakraborty, A. K., Kardar, M. & Shakhnovich, E. I. Thymic selection of T-cell receptors as an extreme value problem. *Phys. Rev. Lett.* **103**, 068103 (2009).
- Chao, D. L., Davenport, M. P., Forrest, S. & Perelson, A. S. The effects of thymic selection on the range of T cell cross-reactivity. *Eur. J. Immunol.* **35**, 3452–3459 (2005).
- Naeher, D. *et al.* A constant affinity threshold for T cell tolerance. *J. Exp. Med.* **204**, 2553–2559 (2007).

13. Turnbull, E. L. *et al.* HIV-1 epitope-specific CD8⁺ T cell responses strongly associated with delayed disease progression cross-recognize epitope variants efficiently. *J. Immunol.* **176**, 6130–6146 (2006).
14. Gillespie, G. M. *et al.* Cross-reactive cytotoxic T lymphocytes against a HIV-1 p24 epitope in slow progressors with B*57. *AIDS* **16**, 961–972 (2002).
15. Yu, X. G. *et al.* Mutually exclusive T-cell receptor induction and differential susceptibility to human immunodeficiency virus type 1 mutational escape associated with a two-amino-acid difference between HLA class I subtypes. *J. Virol.* **81**, 1619–1631 (2007).
16. Moon, J. J. *et al.* Naive CD4⁺ T cell frequency varies for different epitopes and predicts repertoire diversity and response magnitude. *Immunity* **27**, 203–213 (2007).
17. Altfeld, M. *et al.* HLA alleles associated with delayed progression to AIDS contribute strongly to the initial CD8⁺ T cell response against HIV-1. *PLoS Med.* **3**, e403 (2006).
18. Althaus, C. L. & De Boer, R. J. Dynamics of immune escape during HIV/SIV infection. *PLoS Comput. Biol.* **4**, e1000103 (2008).
19. Nowak, M. A. *et al.* Antigenic oscillations and shifting immunodominance in HIV-1 infections. *Nature* **375**, 606–611 (1995).
20. Wodarz, D. & Thomsen, A. R. Effect of the CTL proliferation program on virus dynamics. *Int. Immunol.* **17**, 1269–1276 (2005).
21. Cao, J. H., McNevin, J., Malhotra, U. & McElrath, M. J. Evolution of CD8⁺ T cell immunity and viral escape following acute HIV-1 infection. *J. Immunol.* **171**, 3837–3846 (2003).
22. Price, D. A. *et al.* Public clonotype usage identifies protective Gag-specific CD8⁺ T cell responses in SIV infection. *J. Exp. Med.* **206**, 923–936 (2009).
23. Kiepiela, P. *et al.* Dominant influence of HLA-B in mediating the potential co-evolution of HIV and HLA. *Nature* **432**, 769–775 (2004).
24. Thio, C. L. *et al.* HLA-Cw*04 and hepatitis C virus persistence. *J. Virol.* **76**, 4792–4797 (2002).
25. McKiernan, S. M. *et al.* Distinct MHC class I and II alleles are associated with hepatitis C viral clearance, originating from a single source. *Hepatology* **40**, 108–114 (2004).
26. Bhalerao, J. & Bowcock, A. M. The genetics of psoriasis: a complex disorder of the skin and immune system. *Hum. Mol. Genet.* **7**, 1537–1545 (1998).
27. Chessman, D. *et al.* Human leukocyte antigen class I-restricted activation of CD8⁺ T cells provides the immunogenetic basis of a systemic drug hypersensitivity. *Immunity* **28**, 822–832 (2008).
28. López de Castro, J. A. HLA-B27 and the pathogenesis of spondyloarthropathies. *Immunol. Lett.* **108**, 27–33 (2007).
29. Bowness, P. HLA B27 in health and disease: a double-edged sword? *Rheumatology* **41**, 857–868 (2002).
30. Streeck, H. *et al.* Human immunodeficiency virus type 1-specific CD8⁺ T-cell responses during primary infection are major determinants of the viral set point and loss of CD4⁺ T cells. *J. Virol.* **83**, 7641–7648 (2009).

Supplementary Information is linked to the online version of the paper at www.nature.com/nature.

Acknowledgements Financial support was provided by the Mark and Lisa Schwartz Foundation, the National Institutes of Health (NIH) Director's Pioneer award (A.K.C.), Philip T and Susan M Ragon Foundation, Jane Coffin Childs Foundation (E.L.R.), the Bill and Melinda Gates Foundation, and the NIAID (B.D.W., T.M.A. and M.A.). This project has been funded in whole or in part with federal funds from the National Cancer Institute, NIH, under contract no. HHSN261200800001E. The content of this publication does not necessarily reflect the views or policies of the Department of Health and Human Services, nor does mention of trade names, commercial products, or organizations imply endorsement by the US Government. This research was supported in part by the Intramural Research Program of the NIH, National Cancer Institute, Center for Cancer Research.

Author Contributions A.K. and E.L.R. contributed equally to this work. A.K.C. and B.D.W. initiated the project. A.K., E.L.R. and A.K.C. developed the computational models. A.K., E.L.R., A.K.C. and B.D.W. analysed computational results. Y.Q., F.P., M.C., S.G.D. and B.D.W. collected and analysed the data from cohorts of HIV-infected people. A.K., E.L.R., T.M.A., M.A., M.C., B.D.W. and A.K.C. contributed to the writing of the manuscript.

Author Information Reprints and permissions information is available at www.nature.com/reprints. The authors declare no competing financial interests. Correspondence and requests for materials should be addressed to A.K.C. (arupc@mit.edu) or B.D.W. (bwalker@partners.org).

METHODS

HLA-peptide binding predictions. There are at present several HLA-peptide binding prediction methods. The performance of these algorithms to identify new epitopes has recently been benchmarked against experimental data³¹. In general, artificial neural networks (ANN)³² and the stabilized matrix method (SMM)³³ were found to be superior to other methods³¹. We used ANN and the SMM (versions 2009-09-01 and 2007-12-27) prediction tools provided by the IEDB⁴. Accuracy of prediction tools was tested against experimental data downloaded from the IEDB in September 2009 (Supplementary Fig. 1, Supplementary Table 1 and Supplementary Notes 1). These experimental data were obtained by two methods: competition assays, in which purified MHC and radioactive labelling are used; and association studies, in which purified MHC and fluorescence labelling are used. Data obtained from the two methods show significant correlations of measured binding affinities (as measured by half-maximum inhibitory concentration (IC₅₀) and half-maximum effective concentration (EC₅₀))⁵. Prediction tools were tested against experimental data for accuracy of classifying peptides into binders (IC₅₀ < 500 nM) and non-binders (IC₅₀ ≥ 500 nM); the chosen thresholds are commonly accepted values⁵. We also tested how well these tools predict absolute measured affinity values, not just classification of binders and non-binders, which is dependent on the chosen thresholds. The accuracy of the prediction tools thus determined are summarized in Supplementary Table 1 and Supplementary Fig. 1. We excluded all HLA and Mamu alleles for which there was not enough experimental data (at least 50 binders and 50 non-binders) or prediction tools were not sufficiently accurate (Supplementary Notes 1). For each HLA and Mamu allele, the most accurate prediction tool was used to predict the fraction of unique peptides derived from the human and macaque proteome (Homo_sapiens.GRCh37.55.pep.all.fa and Macaca_mulatta.MMUL_1.56.pep.all.fa obtained from Ensembl⁶) that can bind to that allele. We focused only on the binding abilities of peptides of 9 amino acids to HLA molecules, because there is not enough experimental data available for the binding affinities of peptides of 8, 10 and 11 amino acids to HLA-B*5701 and the other relevant HLA-B alleles that emerged from our analyses (HLA-B*2705, HLA-B*0702 and HLA-B*3501).

Thymic selection model and antigen recognition. The TCR contact residues of peptides and the peptide contact residues of TCRs are represented as strings of sites of length N . One-million sequences of TCR peptide contact residues were subject to development in a thymus containing M self peptides with TCR contact residues generated according to their frequency of occurrence in the human proteome. A particular TCR with the sequence of peptide contact residues \vec{t} successfully matures in the thymus if it avoids negative selection with all self peptides ($E_{\text{int}} > E_n$) and is positively selected by at least one self peptide ($E_{\text{int}} < E_p$). Interaction free energy between sequences of TCR and peptide contacts, \vec{t} and \vec{s} , is, respectively:

$$E_{\text{int}}(\vec{t}, \vec{s}) = E_c + \sum_{i=1}^N J(t_i, s_i) \quad (1)$$

where E_c represents an interaction between a TCR and an HLA molecule, and J is an empirically determined statistical potential between interacting amino acids on a TCR and a peptide. Antigenic peptides are recognized by a mature TCR if binding is stronger than the threshold for recognition ($E_{\text{int}} < E_r$). The statistical potentials do not necessarily provide quantitatively accurate values of the interaction free energies. However, theoretical analyses and computational results^{9,10} show that the following qualitative result is true regardless of the choice of the statistical potentials: the smaller the diversity of self peptides presented in the thymus, the greater the cross-reactivity of the mature T-cell repertoire that develops therein. More details of the model and the insensitivity of our results to parameter variations (for example, qualitative results do not depend on the choice of J or E_c (as long as E_c is not too small or large)) are described in Supplementary Information (Supplementary Figs 2 and 3) and elsewhere^{9,10}. The parameters used for the results in the main text are: $N = 5$; $E_n - E_c = -21 k_B T$; $E_p - E_n = 2.5 k_B T$; $E_r = E_n$ and Miyazawa-Jernigan statistical potential J^{34} . Numbers of self peptides presented in the thymus, M , were varied to represent different HLA alleles.

Host-pathogen interaction dynamics. We constructed a small model of HIV with distinct epitopes and sequence diversity, based in part on models developed previously^{18,19}. The virus is modelled as displaying L epitopes, each consisting of M amino acid residues that may be of N types. Different viral strains arise through point mutations at the amino-acid sites, giving $(N^M)^L$ distinct strains. The number of different pMHC types is $L \times N^M$, because peptide sequences at epitope positions $1 \dots L$ are considered to be distinct. The system of ordinary differential equations corresponding to the model in Fig. 2 and based on previous work²⁰ is as follows:

$$\frac{dV_n}{dt} = k_v^n I_n - k_c V_n + k_m \sum_{n,m} (V_m - V_n) \quad (2)$$

$$\frac{dI^t}{dt} = k_b - k_d I^t - k_t I^t \sum_n V_n \quad (3)$$

$$\frac{dI_n}{dt} = k_t V_n I^t - k_d' I_n - \sum_i \sum_j \sigma_{ij} k_k P_{n,j} T_i^* \quad (4)$$

$$\frac{dP_{n,j}}{dt} = k_s I_n - k_o P_{n,j} - \frac{dI_n^{(\text{kill})}}{dt} \frac{P_{n,j}}{I_n} \quad (5)$$

$$\frac{dP_{n,j}^{\text{APC}}}{dt} = k_s' I_n - k_o' P_{n,j}^{\text{APC}} \quad (6)$$

$$\frac{dT_i}{dt} = -k_a T_i \sum_{n,j} \sigma_{ij} P_{n,j}^{\text{APC}} \quad (7)$$

$$\frac{dT_i^0}{dt} = -k_p T_i^0 + k_a T_i \sum_{n,j} \sigma_{ij} P_{n,j}^{\text{APC}} + k_{ra} M_i \sum_{n,j} \sigma_{ij} P_{n,j}^{\text{APC}} \quad (8)$$

$$\frac{dT_i^m}{dt} = 2k_p T_i^{(m-1)} - k_p T_i^m \quad (9)$$

$$\frac{dT_i^*}{dt} = 2k_p T_i^{(D-1)} - k_{dt} T_i^* - k_m T_i^* \quad (10)$$

$$\frac{dM_i}{dt} = k_m T_i^* - k_{dm} M_i - k_{ra} M_i \sum_{n,j} \sigma_{ij} P_{n,j}^{\text{APC}} \quad (11)$$

Target CD4⁺ T cells, I^t , are infected by free virus particles, where V_n denotes virions of strain n . I_n denotes CD4⁺ T cells infected by virus of strain n , $P_{n,j}$ is a pMHC complex of peptide j derived from viral strain n , displayed on the surface of the infected cell, $P_{n,j}^{\text{APC}}$ is a pMHC displayed by APCs and T_i is a naive CD8⁺ T cell of clonotype i . Activated T cells undergo D rounds of cell division before becoming effector CTLs; T_i^0 is an activated CD8⁺ T cell of type i that has not yet begun dividing and T_i^m are the dividing cells, where m runs from 1 to $D-1$. Effector CTLs, T_i^* , differentiate into memory CD8⁺ T cells, M_i , which are activated upon re-exposure to pMHC.

If T-cell clone i recognizes pMHC j , σ_{ij} is 1, and 0 otherwise. In equation (2), $\sum_{n,m}$ denotes the sum over viral strains m that are Hamming distance 1 away from strain n . That is, only point mutations are allowed. The third term of equation (5) ensures that if an infected cell is killed, the pMHC bound on its surface must also disappear; $\frac{dI_n^{(\text{kill})}}{dt}$ denotes the third term of equation (4), which describes killing of an infected cell by CTLs that recognize pMHC on its surface. Simulations were performed using ode45 and ode15s solvers in MATLAB. A further dynamic model, which does not incorporate target cell limitation and allows unlimited expansion of activated CTLs, was also developed to show robustness of our results to model assumptions. It is discussed in the Supplementary Information (Supplementary Figs 7–12).

Rate constants used in the models are given in Supplementary Tables 3–4, and are in keeping with values reported in the literature. We assume a concentration of 10^6 CD4⁺ T cells per ml blood before infection, with 1% of these cells activated and thus initial targets for HIV infection^{35,36}. The initial conditions of infection in the simulations were one infected CD4⁺ T cell per ml of plasma and a naive-CD8⁺ repertoire size of one cell per ml of each clonotype. We assume that the number of epitopes, length of each epitope, and number of amino acids (L, M, N) are all 2, giving 8 pMHC types and 16 possible viral strains. The number of CD8⁺ clonotypes was chosen to be 20.

The interplay between antigen and immune receptor diversity is captured in this model through variability in σ_{ij} and viral fitness. Different fitness levels for different strains of the virus are modelled by randomly selecting k_v^n , the virus proliferation rate, for each strain from a uniform distribution between 0 and 2,000 per day^{18,37}, with the assumption that the infecting strain has the maximum fitness. The matrix σ_{ij} encodes the ability of T cells to recognize pMHCs. We generate σ_{ij} in such a way as to mimic the results of the thymic selection model (Fig. 1b), to investigate the effects of those predictions on host-pathogen dynamics. That is, we assume that T-cell repertoires restricted by different HLA types differ in the interaction free energies of their TCR-pMHC contacts, and generate σ_{ij} accordingly using a type of random-energy-like model (Supplementary Fig. 6). The interaction free energy between a T cell and an epitope is given by $\sum_a J(i, j_a)$, where $J(i, j_a)$ is

the interaction free energy between T cell of clonotype i and residue a on epitope j . Similar to the models used for thymic selection, the total interaction free energy is taken to be the sum of the individual residue interactions and recognition is said to occur when it exceeds a recognition threshold (in the dynamic model, T-cell sequences are not specified explicitly). $J(i, j_a)$ is a random variable chosen from a uniform distribution, and the width of the distribution determines the probability that the summed interaction energy falls above the threshold, and thus the probability that a peptide is recognized by a given T cell. Repertoires generated in this way approach a Gaussian distribution of interaction energies, and the distribution shifts and thus cross-reactivity increases when the uniform distribution from which $J(i, j_a)$ is selected is wider. Generating $\sigma_{i,j}$ in this way allows us to describe variable cross-reactivities of the T-cell repertoire (both intra- and inter-epitope), and also accounts for correlated interaction energies and thus recognition probabilities of similar peptide sequences.

HLA-allele association with ability to control HIV. SAS 9.1 (SAS Institute) was used for data management and statistical analyses. Odds ratios and 95% confidence intervals were determined using PROC LOGISTIC in a comparison of HIV controllers (those individuals who maintained viral loads of less than 2,000 copies of the virus per ml plasma on three determinations over at least a year of follow-up and, on average, for approximately 15 years³⁸) to HIV non-controllers (those individuals whose viral loads exceeded 10,000 copies of the virus per ml plasma). To eliminate the confounding effects of B*0702, B*3501, B*2705 and B*5701, alleles strongly associated with progression or control, these factors were used as covariates in the logistic regression model for the analysis of all other HLA class I types³⁹. All ethnic groups were included in the analyses shown (European, African-American and others) and we adjusted for ethnicity in the

logistical regression model. All P values were corrected for multiple tests using the Bonferroni correction, a stringent and commonly used approach for multiple comparisons⁴⁰.

31. Peters, B. *et al.* A community resource benchmarking predictions of peptide binding to MHC-I molecules. *PLoS Comput. Biol.* **2**, e65 (2006).
32. Gulukota, K., Sidney, J., Sette, A. & DeLisi, C. Two complementary methods for predicting peptides binding major histocompatibility complex molecules. *J. Mol. Biol.* **267**, 1258–1267 (1997).
33. Peters, B., Tong, W., Sidney, J., Sette, A. & Weng, Z. Examining the independent binding assumption for binding of peptide epitopes to MHC-I molecules. *Bioinformatics* **19**, 1765–1772 (2003).
34. Miyazawa, S. & Jernigan, R. L. Residue-residue potentials with a favorable contact pair term and an unfavorable high packing density term, for simulation and threading. *J. Mol. Biol.* **256**, 623–644 (1996).
35. Sachsenberg, N. *et al.* Turnover of CD4⁺ and CD8⁺ T lymphocytes in HIV-1 infection as measured by Ki-67 antigen. *J. Exp. Med.* **187**, 1295–1303 (1998).
36. Stafford, M. A. *et al.* Modeling plasma virus concentration during primary HIV infection. *J. Theor. Biol.* **203**, 285–301 (2000).
37. Parera, M., Fernandez, G., Clotet, B. & Martinez, M. A. HIV-1 protease catalytic efficiency effects caused by random single amino acid substitutions. *Mol. Biol. Evol.* **24**, 382–387 (2007).
38. Pereyra, F. *et al.* Genetic and immunologic heterogeneity among persons who control HIV infection in the absence of therapy. *J. Infect. Dis.* **197**, 563–571 (2008).
39. Hosmer, D. W., Jovanovic, B. & Lemeshow, S. Best subsets logistic-regression. *Biometrics* **45**, 1265–1270 (1989).
40. Cheverud, J. M. A simple correction for multiple comparisons in interval mapping genome scans. *Heredity* **87**, 52–58 (2001).