

REPORT



# Hiding in plain sight: structure and sequence analysis reveals the importance of the antibody DE loop for antibody-antigen binding

Simon P. Kelow<sup>a,b</sup>, Jared Adolf-Bryfogle<sup>c,d,e</sup>, and Roland L. Dunbrack<sup>a</sup>

<sup>a</sup>Institute for Cancer Research, Fox Chase Cancer Center, Philadelphia, PA, USA; <sup>b</sup>Department of Biochemistry and Molecular Biophysics, University of Pennsylvania, Philadelphia, PA, USA; <sup>c</sup>Protein Design Lab, Institute for Protein Innovation, Boston, MA, USA; <sup>d</sup>Division of Hematology/Oncology, Boston Children's Hospital, Boston, MA, USA; <sup>e</sup>Department of Pediatrics, Harvard Medical School, Boston, MA, USA

## ABSTRACT

Antibody variable domains contain “complementarity-determining regions” (CDRs), the loops that form the antigen binding site. CDRs1-3 are recognized as the canonical CDRs. However, a fourth loop sits adjacent to CDR1 and CDR2 and joins the D and E strands on the antibody v-type fold. This “DE loop” is usually treated as a framework region, even though mutations in the loop affect the conformation of the CDRs and residues in the DE loop occasionally contact antigen. We analyzed the length, structure, and sequence features of all DE loops in the Protein Data Bank (PDB), as well as millions of sequences from HIV-1 infected and naïve patients. We refer to the DE loop as H4 and L4 in the heavy and light chains, respectively. Clustering the backbone conformations of the most common length of L4 (6 residues) reveals four conformations: two  $\kappa$ -only clusters, one  $\lambda$ -only cluster, and one mixed  $\kappa/\lambda$  cluster. Most H4 loops are length-8 and exist primarily in one conformation; a secondary conformation represents a small fraction of H4-8 structures. H4 sequence variability exceeds that of the antibody framework in naïve human high-throughput sequences, and both L4 and H4 sequence variability from  $\lambda$  and heavy germline sequences exceed that of germline framework regions. Finally, we identified dozens of structures in the PDB with insertions in the DE loop, all related to broadly neutralizing HIV-1 antibodies (bNabs), as well as antibody sequences from high-throughput sequencing studies of HIV-infected individuals, illuminating a possible role in humoral immunity to HIV-1.

## ARTICLE HISTORY

Received 16 July 2020  
Revised 8 October 2020  
Accepted 18 October 2020

## KEYWORDS

Antibody therapeutics; antibody structure; structural bioinformatics; antibody complementarity determining regions

## Introduction

Antibodies use three hypervariable loops on each variable domain to bind antigens. These three loops, referred to as complementarity-determining regions (CDRs), were first identified by their high sequence variation relative to the rest of the variable domain sequence.<sup>1</sup> However, there is a fourth loop, structurally adjacent to CDR1 and CDR2 referred to as the DE loop, which joins strands D and E in the immunoglobulin v-type fold (Figure 1).<sup>2,3</sup> In the linear sequence, the DE loop sits between CDRs 2 and 3 and is encoded by V-region gene segments.<sup>4</sup> The DE loop has traditionally been considered part of the antibody framework, so studies addressing the ability of specific DE loop residues to affect antibody binding<sup>5-8</sup> have addressed these residues as framework residues, and not part of a CDR-like loop. However, mutations in the DE loop can affect antigen binding, and in some structures, it directly contacts antigen.


Chothia and Lesk first noted hydrophobic packing interactions of the DE loop with L1, in particular that VL residue 87 (ImMunoGeneTics information system® (IMGT) numbering; Chothia residue 71) packs against L1, and is typically either Phe or Tyr.<sup>5</sup> Foote and Winter demonstrated that some antibodies

lose binding affinity to target antigen upon mutation of this residue from Tyr to Phe, noting that this interaction mediates interaction of L1 with target antigen though a hydrogen bond between Tyr87 and Asn37.<sup>7</sup> Al-Lazikani et al. observed a switch in conformation of CDR L1 of length 11 when Tyr87 changes to Phe87.<sup>8</sup> Tramontano et al. noted that an Arg residue at IMGT VH residue 80 in the heavy chain (Chothia VH residue 71) makes hydrogen bonds to H1 and H2 and packs against side-chain residues of H1 and H2, stabilizing specific H2 conformations and bringing H1 and H2 into closer contact with each other.<sup>6</sup> Several studies since these initial observations have considered various mutations of DE loop residues, with particular focus on VH residue 80, and successfully engineered significant changes in both antibody stability or antibody-antigen affinity.<sup>9-13</sup> However, the effects (or lack thereof) of these mutations are unpredictable, and appear to vary with the germline construct of the antibody.

Previously we demonstrated the importance of the DE loop in redesigning an unstable anti-epidermal growth factor receptor antibody, C10, and its affinity-matured form P2224.<sup>14</sup> The VL region of C10 appeared to be a fusion of  $\lambda$ 3 and  $\lambda$ 1 V-region gene loci, introduced most likely through PCR amplification. We redesigned the antibody framework in an attempt to

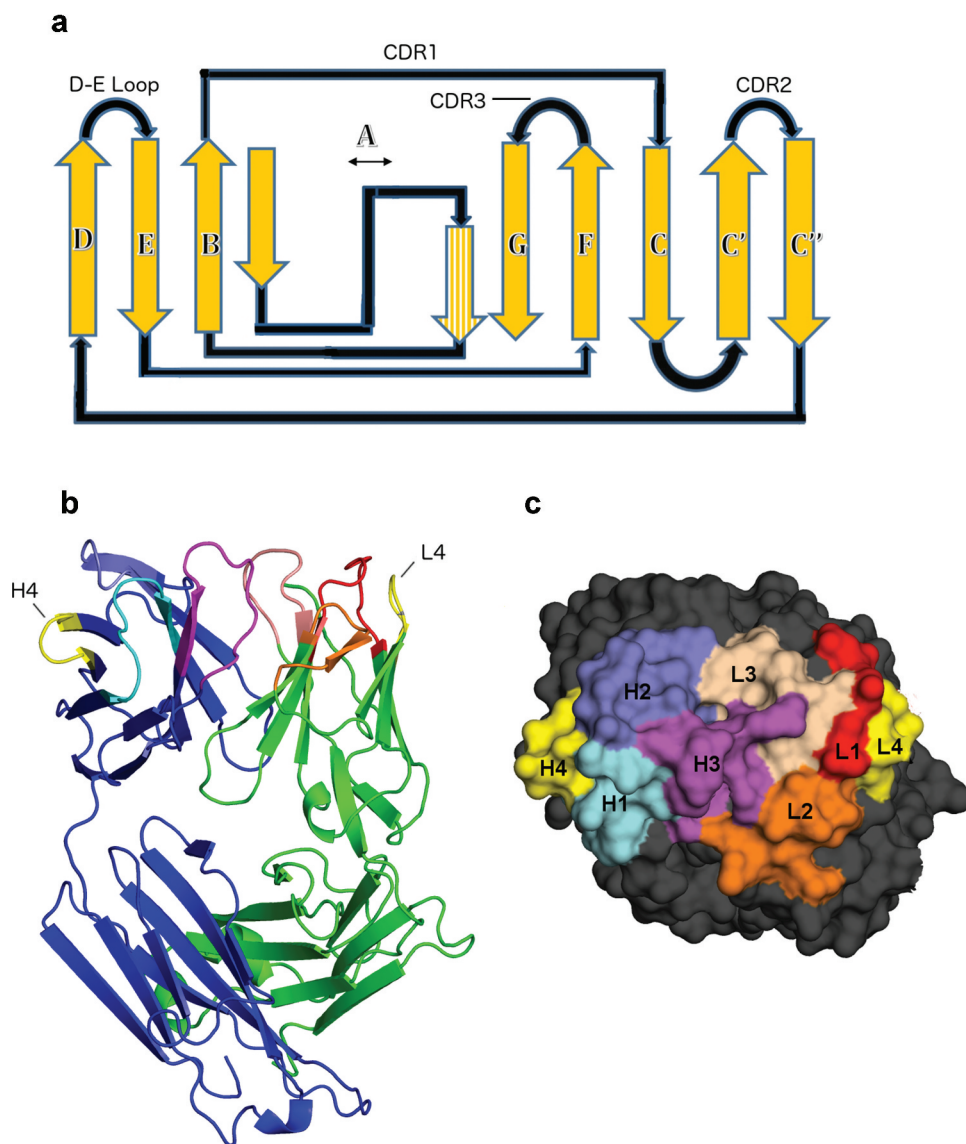
**CONTACT** Roland L. Dunbrack  [Roland.Dunbrack@fcc.edu](mailto:Roland.Dunbrack@fcc.edu)  Institute for Cancer Research, Fox Chase Cancer Center, Philadelphia, PA, USA.

Abbreviations: bNabs: broadly neutralizing antibodies; CDR – complementarity-determining region; DBSCAN: Density-based spatial clustering of applications with noise; DE loop: loop connecting strands D and E of an immunoglobulin variable domain; EDIA: Electron Density for Individual Atoms (program); GP120: envelope glycoprotein in HIV-1 envelope; HIV: human immunodeficiency virus; IMGT: the international ImMunoGeneTics information system® (<http://imgt.org>); PDB: Protein Data Bank

 Supplemental data for this article can be accessed on the [publisher's website](#).

© 2020 The Author(s). Published with license by Taylor & Francis Group, LLC.

This is an Open Access article distributed under the terms of the Creative Commons Attribution-NonCommercial License (<http://creativecommons.org/licenses/by-nc/4.0/>), which permits unrestricted non-commercial use, distribution, and reproduction in any medium, provided the original work is properly cited.



**Figure 1.** Position of the DE loop in antibody structures. (a) V-type fold according to Bork et. al. The DE loop and CDRs are indicated. Strand A forms beta-strand pairing interactions with both strand B and strand G. (b) Example of antibody Fab fragment (light chain in green, heavy chain in blue). (c) Top-down view of antibody combining site. The canonical CDRs and the DE loop are marked in panel C and are represented in the same colors in panel B.

stabilize the antibody and prevent antibody aggregation by grafting the sequences of the  $\lambda$  antibody L1, L2, and L3 CDRs onto a  $\kappa$  framework. We observed that the  $\lambda$  DE loop was different in structure and sequence from a typical  $\kappa$  DE loop in antibodies. Grafting the DE loop along with L1, L2, and L3 from the P2224  $\lambda$  antibody onto a  $\kappa$  framework produced an antibody with significantly increased thermostability, while also retaining P2224's binding affinity. As a control, grafting L1, L2, and L3 while keeping the host  $\kappa$  DE loop sequence produced an antibody with lower stability and significantly reduced affinity.

Here, we analyze the structures and sequences of the DE loops of antibody heavy and light chain variable domains in the Protein Data Bank (PDB), along with a large set of sequences from multiple high-throughput antibody sequencing studies. We first define the DE loop (which we refer to as L4 on the light chain and H4 on the heavy chain) as IMGT residues 80–87 based on the structural variability observed in Ramachandran

maps of residues encompassing the D and E strands of the heavy and light chains in the PDB. With these definitions, we expand on the observations presented by Lehmann et al. by clustering the backbone conformations of L4 and H4 loops in the structures of antibodies to address their structural contribution to antigen binding. The vast majority of L4 loops are of length 6, the exceptions being human  $\lambda$ 5 and  $\lambda$ 6, mouse  $\lambda$ 4– $\lambda$ 8, rat  $\lambda$ 2 and  $\lambda$ 3, and rabbit  $\lambda$ 5 and  $\lambda$ 6 DE loops, which are length 8. All human and mouse germline H4 loops are of length 8. Some rabbit and llama H4 sequences are of length 7.

From a clustering of the conformations of L4 and H4 with validated electron density,<sup>15</sup> we demonstrate that L4 loops of length 6 exist in two dominant conformations, one of which only contains  $\kappa$  antibodies, while the other contains both  $\lambda$  and  $\kappa$  antibodies. There are also two smaller clusters, one of which is  $\kappa$  and one of which is  $\lambda$ . H4 length-8 structures have one primary conformation, as well as a secondary conformation that represents a small fraction of the H4 length-8 structures.

The primary heavy chain H4 cluster and L4 length-8 structures have very similar backbone conformations. In addition to classifying the structures of L4 of lengths 6 and 8 and H4 of lengths 6, 7, and 8, we also calculate all hydrogen bond interactions between the DE loop and the CDRs that influence the conformation of these CDRs. We also correlate the structural features with antibody germline identity as defined by the IMGT database.<sup>4</sup>

Finally, we examine 125 structures in the PDB that have insertions in the heavy or light chain DE loops compared to their germlines. With only one exception, these all turn out to be structures of broadly neutralizing antibodies (bnAbs) of HIV-1 in human patients.<sup>16–52</sup> From these structures, we identify insertions in the DE loops on both the light and heavy chains that contact the antigen gp120, and compare the binding contribution of the DE loops with insertions to the rest of the CDRs in these antibodies. From sequencing studies of HIV-1 infected individuals, we identify insertions and deletions, hypersomatic mutation, and frameshift mutations in and around the DE loop region for the light and heavy chains. These same sequence features are rarely observed in a large set of antibodies from uninfected individuals, and thus may represent a mechanism in humoral immunity to HIV-1.

## Results

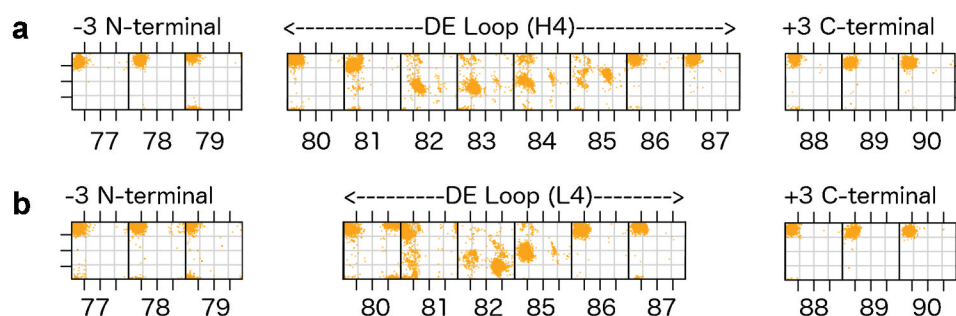
### Clustering of canonical length L4 and H4 structures

To define the regions of structural variability in the vicinity of both H4 and L4, we plotted the  $\phi$  and  $\psi$  dihedrals of the D and E strands and the residues in between for all heavy and light chains of antibodies in the PDB with the most common L4 or

H4 lengths of 6 and 8 residues, respectively (Figure 2). We found that IMGT residues 77–79 and 86–90 uniformly occupy the beta region of the Ramachandran map, while there is some variability in residues 80–82 and 85 of light chains and residues 81–85 of heavy chains. So that the starting and ending residues are opposite each other in the beta strands, and because residues 80 and 87 contact CDR1 as noted above, we define the DE loop as IMGT residues 80–87. Kabat and Chothia number the H4 region as residues 71–78 and L4 loops of length 6 as residues 66–71. L4 loops of length 8 would require insertion codes, such as 68A, 68B. In the rest of this report, we number the residues in the DE loops from 1 to N for DE loops of length N, such that L4 loops of length 6 are numbered 1–6, and L4 and H4 loops of length 8 are numbered 1–8. A mapping of our residue numbering to those of IMGT, Kabat, Chothia, and Honegger-Plückthun is presented in Table 1.

We clustered the structures of L4 loops with germline lengths 6 and 8 and H4 loops of length 6, 7, and 8 using a maximum dihedral angle metric described in Materials and Methods. We used a density-based clustering algorithm, DBSCAN (Density-based spatial clustering of applications with noise)<sup>53</sup> to identify and remove outliers and to identify common conformations within the data. Table 2 provides a summary of the L4 and H4 clusters, specifying their gene, consensus Ramachandran conformation, consensus sequence, number of chains in each cluster, fraction of PDB chains of that length in each cluster, number of unique sequences, and the average  $\phi$  and  $\psi$  dihedral values for each residue in the DE loop for that cluster.

Because structures of low resolution or those with dynamic loops may be solved incorrectly, we used the Electron Density for Individual Atoms (EDIA) program to evaluate the electron



**Figure 2.** Ramachandran plots for part of the D strand, the DE loop, and part of the E strand (IMGT residues 77–90) for the most common DE loop lengths. (a).  $\phi$  (x-axis) and  $\psi$  (y-axis) for residues in length-8 H4 loops, and the 3 anchor residues before and after the loop. IMGT residue numbers are provided at the bottom of each panel. (b).  $\phi, \psi$  plots for residues in length-6 L4 loops, and the 3 anchor residues before and after the loop

**Table 1.** Map between various numbering schemes within H4 and L4 loops.

IMGT	AHo	Position in H4-6	Position in H4-7	Position in H4-8	Chothia/Kabat Heavy Chain	Position in L4-6	Position in L4-8	Chothia/Kabat Light Chain
80	82	1	1	1	71	1	1	66
81	83	2	2	2	72	2	2	67
82	84	3	3	3	73	3	3	68
83	85		4	4	74		4	68A
84	86			5	75		5	68B
85	87	4	5	6	76	4	6	69
86	88	5	6	7	77	5	7	70
87	89	6	7	8	78	6	8	71

Mapping from various antibody numbering schemes to the numbering scheme used (1 to N where N is the length of the CDR loop considered). AHo indicates the Honegger-Plückthun numbering scheme

Table 2. DE loop canonical families.

Gene	Cluster	Rama. string	Common sequence	# PDB chains (all)	# PDB chains ( $\geq 0.75$ EDIA)	% chains (all)	% chains ( $\geq 0.75$ EDIA)	Unique seqs (all)	Unique seqs ( $\geq 0.75$ EDIA)	1	2	3	4	5	6	7	8
k	L4-6-1	EBEABB	GSGTDF	3,854	1,393	64.5	66.6	77	48	120, 170	-164, 158	73, -100	-118, -11	-119, 122	-132, 153		
$\lambda/\kappa^*$	L4-6-2	BBEABB	KSGTTA	1,377	562	23.0	26.9	81	67	-141, 136	-144, 114	64, -120	-100, 10	-120, 136	-115, 148		
k	L4-6-3	BBAABB	GSGTDF	95	48	1.6	2.3	7	4	163, 168	-90, -144	-88, -23	-128, -13	-125, 125	-131, 151		
$\lambda$	L4-6-4	BBLBB	LIGGKA	31	14	0.5	0.7	6	5	-110, 130	-138, 123	51, 48	70, 12	-124, 156	-86, 145		
$\lambda/\kappa$	noise	-	-	622	76	10.4	3.5	74	29	-	-	-	-	-	-		
$\lambda 5/6$	L4-8-1	BBAALBB	IDSSNSA	90	40	100.0	100.0	8	6	-128, 135	-120, 100	-65, -31	-66, -34	-101, 2	54, 44	-134, 156	-111, 151
H	H4-6-1	BBAABB	RTSTTV	35	21	76.1	100.0	4	3	-134, 152	-119, -172	-64, -34	-124, 0	-138, 158	-127, 133		
H	noise	-	-	11	0	23.9	0.0	5	0	-	-	-	-	-	-		
H	H4-7-1	BABAABB	-	37	17	100.0	100.0	7	4	-94, 112	-90, -22	-160, -176	-71, -15	-125, 5	-135, 137	-124, 140	
H	H4-8-1	BBAALBB	RDNSKNTA	6,269	1,953	94.0	96.8	646	333	-143, 149	-122, 110	-65, -30	-67, -34	-102, 2	53, 47	-129, 141	-113, 144
H	H4-8-2	BBAALBB	RDNSKSTA	60	19	0.9	1.1	26	12	-136, 155	-80, 169	-60, -39	-72, -12	64, 28	-104, -11	-135, 137	-115, 145
H	noise	-	-	347	37	5.1	2.1	121	13	-	-	-	-	-	-		

Properties and frequencies of L4 and H4 structural clusters and noise structures for each length of light chain and heavy chain DE loop. The clustering was performed on the entire PDB and on a subset of structures that pass an electron density cutoff ( $EDIA \geq 0.75$ ).  $\phi, \psi$  values (in degrees) are given for each residue in each cluster of L4-6, L4-8, H4-6, H4-7, and H4-8 DE loop length families. Ramachandran map regions are: A = alpha-helix region; B = beta sheet region; E = epsilon region (lower right of Ramachandran map); L = alpha-left region.

\* Cluster L4-6-2 is composed of 75%  $\lambda$  chains and 25%  $\kappa$  chains



density of the backbone atoms in each DE loop in the PDB.<sup>15</sup> EDIA provides an atom-level assessment of electron density by integrating the 2fo-*fc* electron density in a sphere centered on each atomic coordinate. We repeated the clustering on all structures that pass a 0.75 EDIA cutoff for the backbone carbonyl oxygen atoms of the DE loop. The number of structures and unique sequences for each cluster is also provided in Table 2 for the 0.75 EDIA data set, representing structures that are well-solved in the region of the DE loop (see Methods). The average EDIA score for atoms in 2.8 Å structures is 0.8, so most loops with resolution worse than 2.8 Å are removed from the EDIA-filtered data set. The EDIA set is about one third the size of the unfiltered data set, but contains about half the number of unique sequences.

Across the four L4-6 clusters, the changes in backbone conformation may be viewed structurally as a hinge motion away from the variable domain of the antibody, with L4-6-1 being closest to the domain, followed by L4-6-3 and L4-6-2, while L4-6-4 is the farthest away from the domain. Figure 4a shows representative structures of L4-6-1, L4-6-2, L4-6-3, and L4-6-4 DE loops superposed by aligning the stems of the DE loop (−3 C-terminal, +3 N-terminal).

The primary difference between the two biggest clusters, L4-6-1 and L4-6-2, is the amino-acid identity and Ramachandran conformation of the first residue. In the germlines of all human  $\kappa$  light chains and nearly all mouse  $\kappa$  light chains, the first residue of the DE loop is glycine. In human and mouse  $\lambda$  germlines, the first residue is (in order of most common to least common): Lys, Ser, Ile, Asp, Leu, Arg, or Thr in 81 of 84 human and mouse IMGT alleles and Gly in only 3 human alleles of one germline V region (all IGLV9-49, not represented in the PDB). In L4-6-1, the first residue is in an epsilon conformation ( $\phi = 119.8^\circ$ ,  $\psi = 170.2^\circ$ ), consistent with a Gly in  $\kappa$  antibodies. In L4-6-2, the first residue is in a beta conformation ( $\phi = -144.9^\circ$ ,  $\psi = 137.2^\circ$ ), consistent with the  $\lambda$  non-Gly residues. 252 out of 333  $\kappa$  structures in L4-6-2 (75%) contain somatic mutations at the first residue position from Gly to Arg, Ala, Glu, and Gln in decreasing order of frequency. In the clustering of the 0.75 EDIA data set, 86% of the  $\kappa$  structures in L4-6-2 have somatic mutations from Gly to some other amino acid at position 1. Several structures in L4-6-2 have germline-encoded hydrophobic residues (Leu, Ile) at the first two residues of the DE loop in the L4-6-2 sequences (e.g. mouse IGLV1\*01, IGLV2\*01; human IGLV7-43\*01, IGLV7-46\*01).

L4-6-3 is an all- $\kappa$  cluster that differs in conformation from all- $\kappa$  L4-6-1 at positions 2 and 3, such that L4-6-1 has average ( $\phi_2, \psi_2 = -164.1^\circ, 157.5^\circ$ ;  $\phi_3, \psi_3 = 73.6^\circ, -100.0^\circ$ ) and L4-6-3 has average ( $\phi_2, \psi_2 = -90.4^\circ, -137.0^\circ$ ;  $\phi_3, \psi_3 = -93.6^\circ, -26.0^\circ$ ). In the L4-6-3 structures, 18 of 69 (26%) chains (from 3 PDB entries) have somatic mutations at position 3 from Gly to Glu, which accounts for the change in residue 3 from an epsilon conformation to an alpha conformation (with a compensating change at residue 2). The remaining structures have germline sequences, including one structure with a germline Arg residue at residue 3.

L4-6-4 is an all- $\lambda$  cluster that differs from the L4-6-2  $\lambda/\kappa$  cluster at residue positions 3 and 4. This cluster has residues with the left-handed conformation at positions 3 and 4 ( $\phi_3, \psi_3 = 51^\circ, 48^\circ$ ;  $\phi_4, \psi_4 = 70^\circ, 12^\circ$ ), which is not seen in any of the other

clusters. Consistent with these conformations, a majority of the loops have glycine residues at these positions. Similar to some of the sequences in L4-6-2, the sequences that are part of this cluster start with hydrophobic residues at positions 1 and 2 (Leu or Ile) that come from human IGLV7-43, llama IGLV8-3, and mouse IGLV1 and IGLV2 germlines. Of the 31 structures of L4-6-4, 12 structures have the germline sequence. The remaining sequences are somatic mutations of the third position (G→E, 4 structures), the second position (L→R, L→I, 2 structures), and the fourth position (D→G, 13 structures).

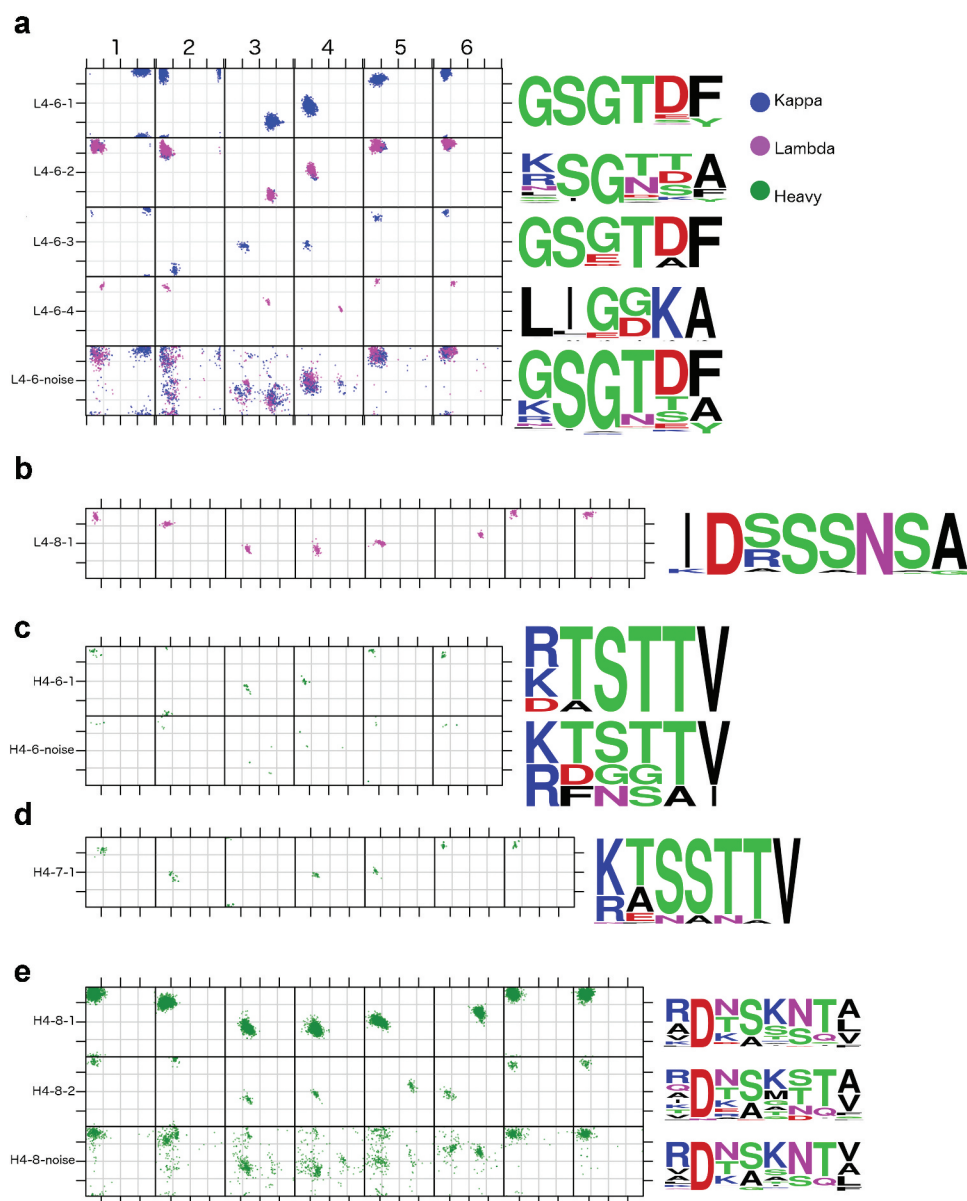
We evaluated both L4-6-3 and L4-6-4 by their electron density by examining how many of their chains and unique sequences are produced from the clustering of the EDIA $\geq$ 0.75 data set. For L4-6-3, four of seven unique sequences remain, and for L4-6-4, five of six unique sequences remain after the EDIA cutoff, indicating that they are robust clusters, and not related to mis-solved residues.

For the 67 length-8 L4 structures in the PDB, which are related to a small number of  $\lambda$  germlines in humans, mice, rats, rabbits, and macaques, we observed a single cluster (Figures 3b, d) representing 7 unique sequences. No structures were placed into noise by the DBSCAN algorithm, indicating a low level of structural variation. Of 49 PDB entries containing light chains with length-8 DE loops, 17 are involved in Bence-Jones homodimers associated with light-chain amyloidosis.<sup>54</sup> After enforcing the EDIA cutoff, 41 of the 67 chains remain, which contain all 6 of the unique sequences found in the clustering with no EDIA cutoff.

For canonical length-8 H4 structures, clustering with DBSCAN produced one large cluster with 6,269 chains, and several small clusters with less than 60 chains each. The sequences of the small clusters were very similar to those in the large cluster. Many of them involve peptide flips from the large cluster and may be incorrectly solved. Peptide flips involve changing the  $\psi$  dihedral of one residue by 180° and the  $\phi$  dihedral of the next residue by 180°.<sup>55</sup> The carbonyl residue of the first residue moves by more than 3 Å in a peptide flip, so the electron density of the oxygen atom of each residue in a loop is diagnostic of mis-solved peptide flips, which are common in protein loops. By clustering the EDIA $\geq$ 0.75 data set, in addition to the large cluster, only one small cluster remained with a substantial number of structures,<sup>19</sup> as well as unique sequences.<sup>12</sup> The others involved peptide flips from the large cluster, and likely are due to incorrectly solved structures.

We therefore chose to name only two H4-8 clusters, H4-8-1 and H4-8-2 (Figures 3e, Figure 4c). The H4-8-1 cluster has 646 unique sequences, exhibiting far greater sequence variation than any of the L4-6 clusters. For H4-8-2 ( $\phi_2, \psi_2 = -80^\circ, -160^\circ$ ) there is a small shift in residue 2 compared to H4-8-1 ( $\phi_2, \psi_2 = -122^\circ, -110^\circ$ ), and a substantial shift at residues 5 and 6, which change the conformation of these residues to LA ( $\phi_5, \psi_5 = 64^\circ, 28^\circ$ ;  $\phi_6, \psi_6 = -104^\circ, -11^\circ$ ), from AL ( $\phi_5, \psi_5 = -102^\circ, 2^\circ$ ;  $\phi_6, \psi_6 = -53^\circ, 47^\circ$ ), respectively. Any other clusters generated in the clustering step had fewer than 10 unique sequences, or they disappeared after enforcing the EDIA cutoff, leading to a suspicion that they are clusters consisting of mis-solved residues at position 5 and 6.

The conformation of length-8 H4 structures is the same conformation as length-8 L4 structures as shown in structural



**Figure 3.** Canonical conformations of L4 and H4.  $\phi, \psi$  plots for each residue in the DE loop for each of the L4 and H4 DE loop clusters from the full data set (no EDIA cutoff) with their respective sequence logos. (a). L4 length-6 loops. (b). L4 length-8 loops. (c). H4 length-6 loops. (d). H4 length-7 loops. (e). H4 length-8 loops

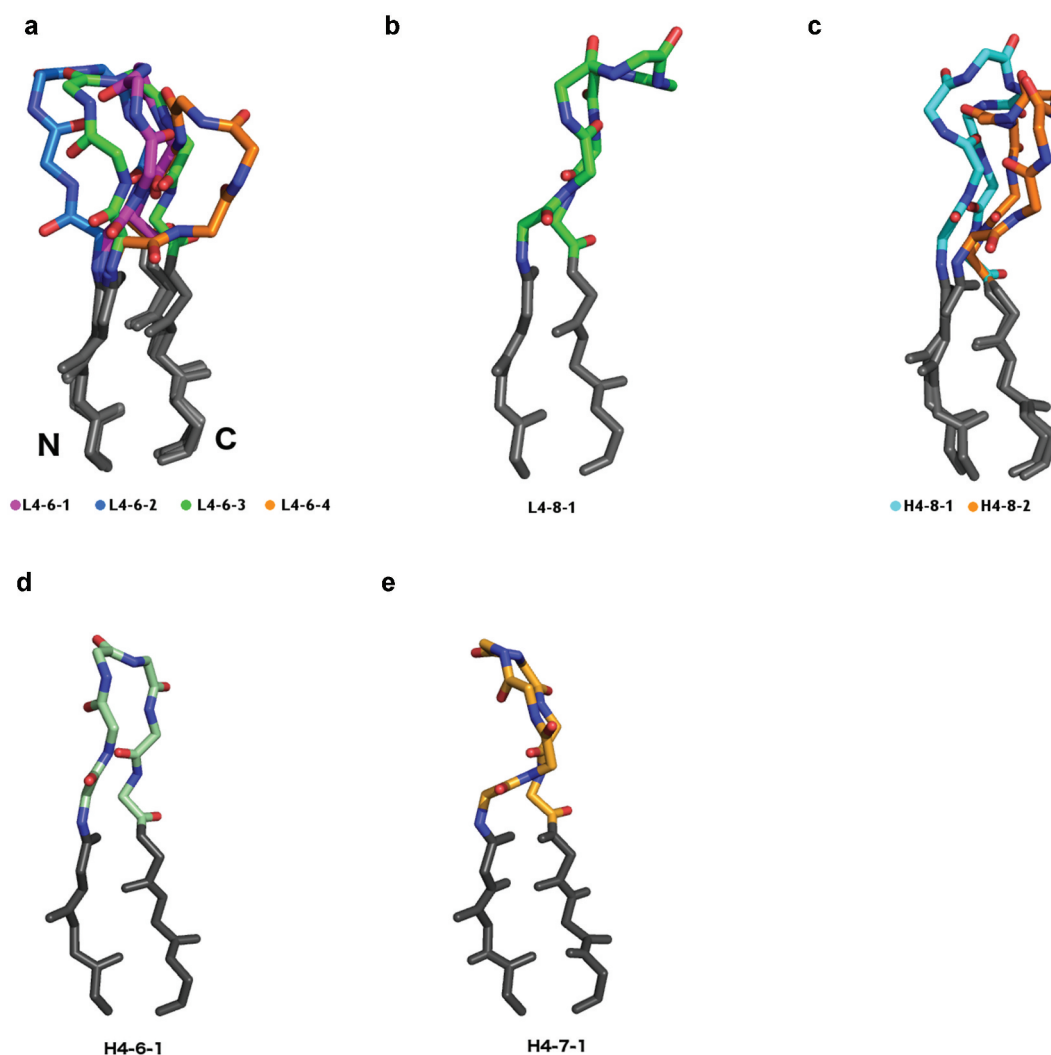
alignment of the two clusters by the CDR4 stem regions (Figure 5a). We also clustered a small number of H4 loops of lengths 6 (from the rabbit IGHV1S69 germline) and 7 (from rabbit IGHV1S45, IGHV1S47, and IGHV1S69 and llama IGHV1S3 germlines) (Table 2, Figure 3c,d), which produced one cluster for each length, designated H4-6-1 and H4-7-1. The H4-6-1 cluster is similar in conformation to L4-6-3, as shown in Figure 5b.

### Interactions between CDR4 and CDRs 1 and 2

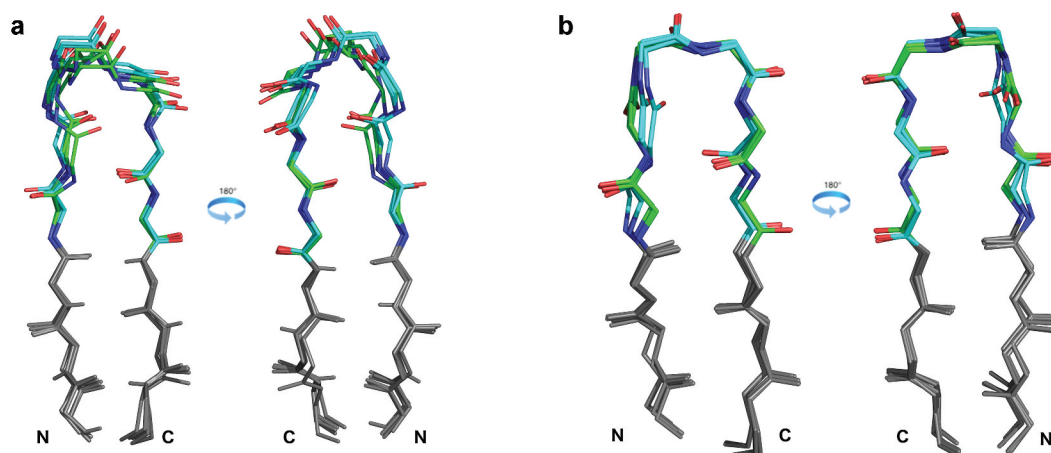
To describe the relationship between various L4 conformations with CDR1 and CDR2 conformations, we first calculated the occurrence of each L4-6 cluster given the various common L1 clusters (Table 3). The  $\kappa$  L1 clusters have more than 80% L4-6-1 DE loops. For most of these, the secondary cluster is L4-6-2, indicating a tendency for residue 1 in the corresponding

germlines to mutate from Gly to another residue type. The L1-16-1 structures prefer L4-6-3 as a secondary cluster, probably because germlines with length 16 more frequently contain residues that are not Gly at position 3 of L4 (11 of 63 human, mouse, and rat germlines in IMGT).

For the  $\lambda$  germlines, all of the L1 clusters except L1-14-1 are 100% L4-6-2. L1-14-1 structures contain 85% L4-6-2 and 15% L4-6-4 DE loops. The L4-6-4 sequence logo (Figure 3a) shows that this cluster comprises sequences that begin with aliphatic amino acids at positions 1 and 2 (Leu or Ile) and have either Lys or Arg at position 5. All germlines with these sequence features have L1 CDRs of length 14 (human IGLV7-43\*01, IGLV7-46\*01/02, IGLV8-61\*01/02; mouse IGLV1\*01/02 and IGLV2\*01; rat IGLV1S1\*01). The charged residue at position 5 does not take part in a hydrogen bond in any of the structures where the Lys or Arg is present with L1-14-1.



**Figure 4.** Structures of all DE loop germline-length clusters. (a) The antibody light-chain DE loop (L4-6) backbone is shown for each cluster where L4-6-1 (PDB 3d9aL, blue) sits closest to the antibody domain, L4-6-3 (1mjuL, green) hinges slightly away and flips the second carbonyl of the DE loop backbone about 180° relative to the other clusters, L4-6-2 (4unuA, magenta) hinges further away from the domain than L4-6-1 and L4-6-3, and L4-6-4 (6frjH, orange) sits the furthest away from the domain. The stems of the DE loop are colored dark gray. (b). Same representation as in (A), but for the sole L4-8-1 (5jpaA, green) cluster. (c). Same representation as in (A), but for the H4-8-1 (5e7bA, cyan) and H4-8-2 (6bliJ, orange) clusters(d). Same representation as in (A), but for the H4-6-1 cluster (6i9iH, green). (e). Same representation as in (A) but for the H4-7-1 cluster (6dbdD).



**Figure 5.** Comparison of H4 and L4 clusters with structural homology. (a). Superposition of several high-resolution heavy chain H4-8-1 structures (cyan, PDB chains: 2x1qA, 4qyoB, 2vxvH), and L4-8-1 structures (green, PDB chains: 1cd0A, 2w01A, 3h0tA) aligned by the stem of CDR4 (colored in gray) show structural similarity between the two clusters. B. Same superposition as in 5A, but for L4-6-3 (green, 1mjuL, 6qnkC, 6mv5L) and H4-6-1 (cyan, 6i8iH, 6c2A, 6banH) clusters.

**Table 3.** Co-occurrence of L1/L4 pairs from structures in the PDB.

L1 cluster	Gene	# chains	L4-6-1	L4-6-2	L4-6-3	L4-6-4	L4-8-1
L1-10-1	κ	185	99.4	0.6	-	-	-
L1-10-2	κ	93	100.0	-	-	-	-
L1-11-1	κ	1559	90.0	9.9	0.1	-	-
L1-11-2	κ	508	91.6	8.2	0.2	-	-
L1-12-1	κ	175	98.1	1.9	-	-	-
L1-12-2	κ	110	98.1	1.9	-	-	-
L1-15-1	κ	274	98.8	-	1.2	-	-
L1-16-1	κ	597	84.1	1.5	14.4	-	-
L1-17-1	κ	277	99.0	0.3	0.7	-	-
L1-11-3	λ	149	-	100.0	-	-	-
L1-12-3	λ	32	-	100.0	-	-	-
L1-13-1	λ	248	-	100.0	-	-	-
L1-13-2	λ	71	-	-	-	-	100.0
L1-14-1	λ	193	-	85.0	-	15.0	-
L1-14-2	λ	176	-	100.0	-	-	-
H1 cluster	Gene	# chains	H4-8-1	H4-8-2	H4-6-1	H4-7-1	
H1-13-1	H	4285	99.1	0.1	0.5	0.3	-
H1-13-2	H	64	96.0	4.0	-	-	-
H1-13-3	H	101	93.0	7.0	-	-	-
H1-13-4	H	183	99.4	0.6	-	-	-
H1-13-5	H	86	93.8	1.2	-	5.0	-
H1-13-6	H	28	100.0	-	-	-	-
H1-13-7	H	58	94.6	5.4	-	-	-
H1-13-10	H	16	81.3	18.7	-	-	-
H1-14-1	H	102	100.0	-	-	-	-
H1-15-1	H	173	97.1	2.9	-	-	-

For each L1 or H1 cluster, the distribution among the L4 or H4 clusters is provided in percent (excluding the noise cluster).

Second, we have calculated all hydrogen bonds between the DE loop and CDR1 or CDR2. Supplementary Table S1 shows all hydrogen bonds calculated between the DE loop and CDR1 or CDR2 for each CDR1 and CDR2 cluster of L1, H1, and H2 (there are no characteristic hydrogen bonds between L4 and L2 with an occupancy over 60%). In the case of hydrogen bonds involving side-chain atoms, the hydrogen bonds are grouped by structures with the same amino acid at the same position within the DE loop.

Hydrogen bonds between the DE loop and CDR1 and CDR2 partition into five broad categories: 1) backbone-backbone hydrogen bonds shared across several CDR1 clusters (Figure 6c); 2) backbone-backbone hydrogen bonds unique to specific DE loop/CDR1 cluster pairings (Figure 6d,e); 3) hydrogen bonds between DE-loop side-chain atoms and backbone atoms at positions shared across several CDR1/CDR2 clusters (Figure 6f); 4) hydrogen bonds between DE loop side-chain atoms and CDR backbone atoms that are specific for some CDR1/2 clusters and lengths (Figure 6a,b,g); and 5) hydrogen bonds between DE loop backbone atoms and CDR1 side-chain atoms that occur in L1 loops longer than 14 residues (Figure 6h,i).

Regardless of DE loop conformation, DE residue 4 in length-6 L4 forms a backbone-backbone hydrogen bond to the backbone nitrogen of the second residue in CDR1 (counting L1 residues immediately after the Cys of the disulfide bond; Figure 6c) for the vast majority of L1 clusters (L1-14-1 excluded). This hydrogen bond is part of the beta sheet containing the C-terminal segment of CDR4 and the N-terminal strand at the beginning of CDR1. On the light chain, most DE loop structures also have a backbone-backbone hydrogen bond

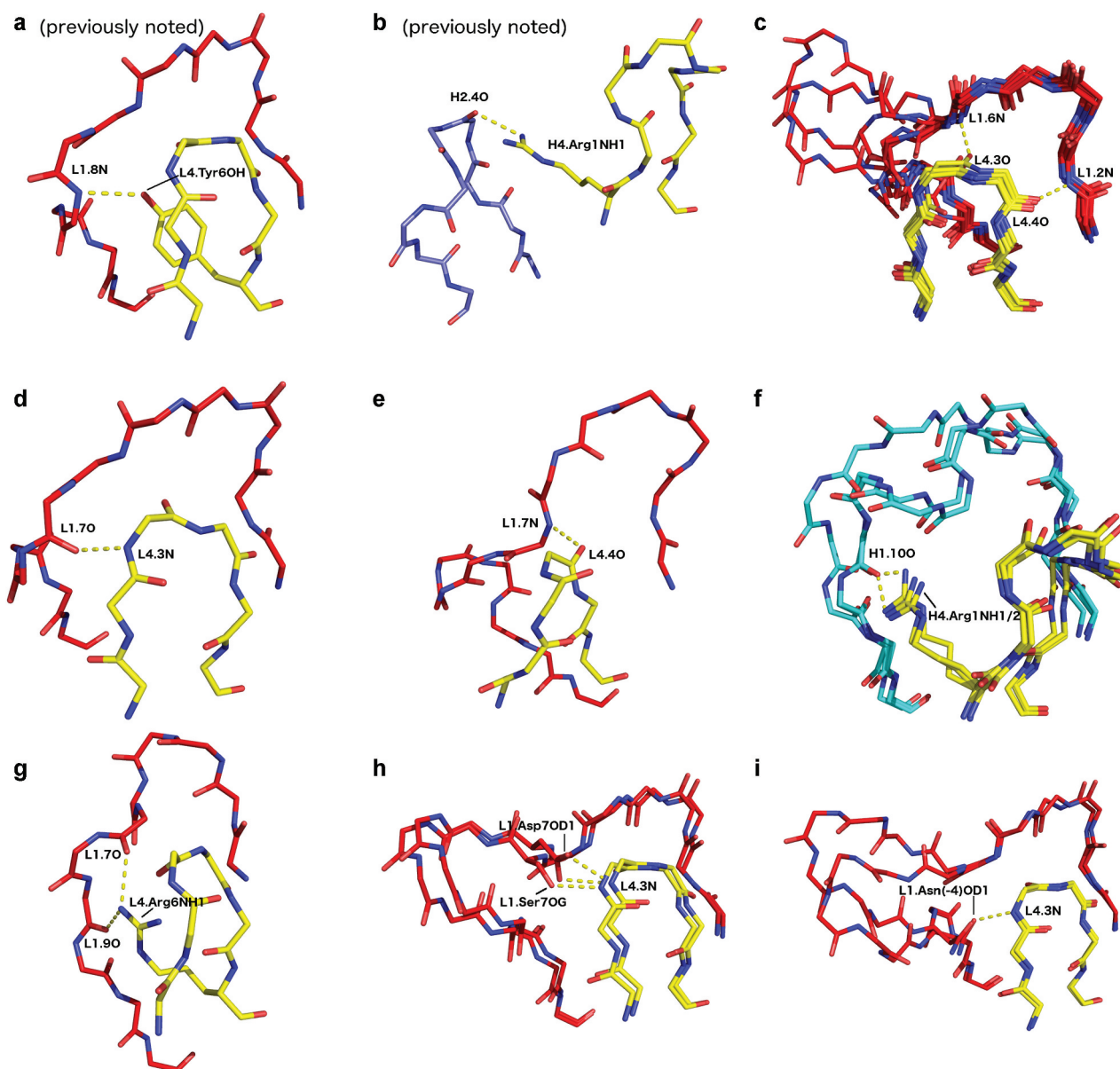
between DE residue 3 and residue 6 of CDR1 (L1-12-3, L1-13-1, L1-14-1, and L1-14-2 excluded; Figure 6c). Structures that have both of these hydrogen bonds have very similar conformations between the residues that are hydrogen bonded, even amongst a diverse set of L1 lengths.

Beyond backbone-backbone hydrogen bonds correlated with the arrangement of the L4 backbone atoms, we note several particular side-chain/backbone hydrogen bonds that occur uniquely with L1 conformations. For example, residue Arg1 hydrogen bonds to the backbone carbonyl of residue 7 of L1-13-1 in 24/26 chains, whereas when the DE residue 1 is Lys, but is still paired with L1-13-1, the hydrogen bond occupancy is only 24% (213 cases).

As noted in previous studies,<sup>5,56</sup> the OH atom of the Tyr6 side chain of the L4 loop forms a hydrogen bond to the backbone nitrogen atom of residue 8 in length-11 L1 CDRs (Figure 6a), flipping its conformation from L1-11-1 (predominantly Phe6) to L1-11-2 (predominantly Tyr6). This hydrogen bond forms in 85% of structures of L1-11-2 with a Tyr residue at position 6 of L4. When this residue is Phe6 instead, this hydrogen bond is lost, and the structure of L1 is L1-11-1, and a hydrogen bond instead forms between the backbone nitrogen of DE residue 3 to the carbonyl of L1-11 residue 7 (Figure 6d). In similar fashion, we note a new hydrogen bond of the side chain of Arg6 in L4-6-2 to the carbonyl backbone oxygen atoms of residues 7 and 9 of L1-12-3 structures, creating a hydrogen bond network. This is an example where the exclusive occurrence of an L1/L4 pair is associated with a unique contact between L1 and L4.

For L1-15-1 and L1-17-1, the carbonyl oxygen of DE residue 1 of L4 is not only hydrogen bonded to a backbone nitrogen





**Figure 6.** Various characteristic hydrogen bonds between the DE loop and CDR1. Hydrogen bonds are labeled CDR4.resnumAtom/CDRn-cluster.resnumAtom (e.g. H4.6O/H1-13-1.2N). If a hydrogen bond is specific to a particular cluster, that is included in the nomenclature. (a). Side-chain/backbone hydrogen bond L4.Tyr6OH/L1-11-2.8N (yellow/red) common in L1-11-2, (b). Side-chain/backbone hydrogen bond H4.Arg1NH1/H2-10-1.3O (yellow/purple) common in H2-10-1. (c). Side-chain/backbone hydrogen bonds L4.3O/L1.6N and L4.4O/L1.2N (both yellow/red) common in L1-10, L1-11, L1-12, L1-13-1, L1-14-2, L1-15-1, L1-16-1, L1-17-1 clusters. (d). Backbone/backbone hydrogen bond L4.3N/L1-11-1.7O common in L1-11-1. (e). Backbone/backbone hydrogen bond L4-6-2.4O/L1-14-1.7N common in L1-14-1. (f). Side-chain/backbone hydrogen bond H4.Arg1NH/H1-13.100 common in H1-13-1,2,3,4. (g). Side-chain/backbone hydrogen bond L4.R6NH/L1-12-3.7O and L4.R6NH/L1-12-3.9O common in L1-12-3. (h). backbone/side-chain hydrogen bonds L4.3N/L1-15-1.57OG and L4.3N/L1-15-1.D7OD common in L1-15-1. (i). L4.3N/L1.N(-4)OD common in L1-16-1 and L1-17-1.

atom in L1, but the backbone nitrogen atom of DE residue 1 is also hydrogen bonded to various side-chain oxygen atoms of residue 7 in L1-15-1 (Asp, Ser, or Thr; [Figure 6h](#)), or residue 14 (Ser or Asn; [Figure 6i](#)) in L1-17-1. Taken together, these results demonstrate that L1/L4 pairs often entail highly specific interactions, facilitated by the L4 cluster-specific arrangement of the L4 backbone and side-chain atoms, which can provide stabilizing hydrogen bonds between L4 and L1.

For H4, in addition to the conserved hydrogen bond involving residue 4 in most H1/H4 pairs for common H1 lengths and clusters (similar to the L4 residue 4 hydrogen bond in [Figure 6c](#)), there are several side-chain/backbone hydrogen bonds that are shared between several H1 clusters and various residues in H4.

Most notably, the Arg1 residue in H4 hydrogen bonds with the backbone carbonyl of residue 10 in the large H1-13-1 cluster using both the NH1 and NH2 atom in the interaction ([Figure 6f](#)). For specific hydrogen bonds, the occupancy of the hydrogen bond depends highly on the H4 residue type. DE residue Asn6 uses its side-chain oxygen atom as well as its side-chain NH2 group to form side-chain/backbone hydrogen bonds between residue 2, residue 5, and residue 7 of H1, stabilizing the H1-14-1 conformation with a hydrogen bond network.

The common Arg1 side-chain of H4 often forms hydrogen bonds with H1 and H2 ([Figure 6b](#)). The most common hydrogen bond of Arg1 is to residue 10 of over 1400 H1-13 structures (in clusters H1-13-1, 2, 3, and 4). In our clustering of the CDR



loop conformations, we noted the common presence of H4-Arg1 in H2-10-2 structures and its relative absence in H2-10-1 structures,<sup>56</sup> implicating a hydrogen bond of Arg1 with the backbone carbonyl of residue 3 of H2-10-2. While the Arg1 side chain is in contact with this carbonyl atom in many structures, the hydrogen bond geometry is poor and Rosetta does not identify it as a hydrogen bond. Its interaction with H2-10-2 is primarily hydrophobic, as noted by Tramontano et al.<sup>6</sup>

For H2-10-1, most H4 loops have an Ala residue at position 1. However, when position 1 in H4 is Arg, 70% of structures have a hydrogen bond from the NH1/NH2/NE atom of Arg to the backbone carbonyl of the residue at position 4 in H2-10 (354 cases with Arg1 when the H2 conformation is H2-10-1). When this residue is Gln, 59% of structures (15/27 chains with H2-10-1 and H4-8-1 with Gln at position 1) contain a hydrogen bond between the NE2 atom of Gln and the backbone carbonyl oxygen of residue 4 of H2. When this residue is instead Lys1, the occupancy of this hydrogen bond is only 42% (17/40 chains with H2-10-1 and H4-8-1 with Lys at position 1). In all other H2-10-1 structures, there is no polar residue at position 1 with a side-chain available to make a hydrogen bond to H2-10. This indicates that a hydrogen bond from position 1 in H4 is not a strong association to the presence of the H2-10-1 conformation, but structures that have this hydrogen bond may have better stability, especially when that hydrogen bond is made from Arg1.

When the conformation of H2-10 is H2-10-2, the Arg1 of H4 rarely forms hydrogen bonds with the carbonyl oxygen at position 4 of H2-10-2 (less than 3% in all 1,180 structures with Arg1 and H2-10-2). Instead, the ND2 atom of Asn3 of H4, which is fully conserved in structures with H2-10-2, forms hydrogen bonds with the same backbone carbonyl oxygen at position 4 of H2 (73% of 1180 structures with Asn at position 3 in H4 and H2-10-2). This points to Asn3 in H4 as a major indicator of the H2-10-2 conformation.

### **Analysis of the sequence variability in DE loops arising from somatically mutated and germline sequences**

From a set of ~2.5 million sequences of naïve human antibodies,<sup>57–62</sup> we calculated the sequence entropy in four of the most prevalent human germlines for the heavy,  $\kappa$ , and  $\lambda$  genes in the data set (Figure 7a), as well as the entropy of human germline sequences of the same length (Figure 7b). As other studies have noted,<sup>63</sup> variability of both framework and CDR residues depends highly upon germline.

For H4 sequence variability, we find that in cases of somatic mutation of any one particular germline, the average sequence entropy of H4 for each of the four germlines exceeds the average sequence entropy for FR1, FR2, and FR3 of the same germline antibodies (Table 4). However, these DE loop residues are less variable than H1 or H2 residues within the same germlines.

For  $\kappa$  antibodies, residue 2 has a higher sequence entropy than about 80% of the other framework residues (Figure 7a). Comparing 28 germline sequences for human  $\kappa$  antibodies (Figure 7b), we observe three highly variable residues (DE residue 2, 5, and 6), and 3 completely conserved residues (IMGT residues Gly1, Gly3 and Thr4). In the variable residues

of L4, the entropy is comparable to the most variable framework residues in germlines, the average entropy does not compare to the average entropy of L1 or L2, and does not exceed the average entropy of FR1, FR2, or FR3 (Table 4).

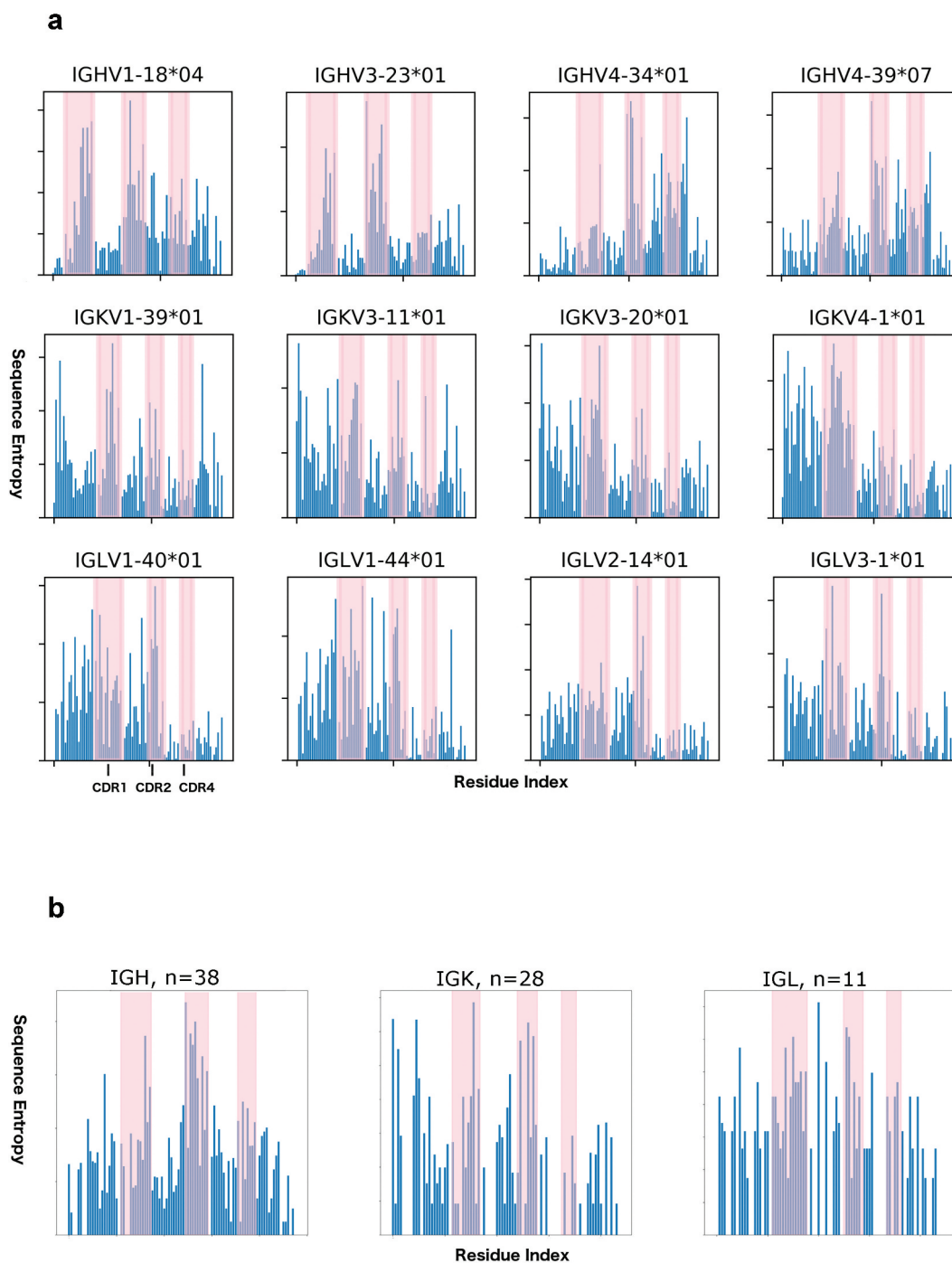
Within each  $\lambda$  germline, L4 sequences are much less somatically mutated than in  $\kappa$  antibody sequences (Figure 7a). The amount of sequence variability due to somatic mutation is less than even the most variable framework residues, and does not compare to sequence variability in L1 or L2. However, looking at 11 germline sequences, average sequence entropy in  $\lambda$  L4 sequences exceeds that of FR1, FR2, and FR3, but is less than that of L1 and L2 (Table 4). Sequence variability at positions 1, 4, and 5 is comparable to both L1 and L2 (Figure 7b). This indicates that sequence variability in  $\lambda$  L4 relates primarily to germline sequence differences, and not somatic mutation.

### **Non-canonical L4 and H4 length in HIV-1 bnAbs**

All known mammalian VH germlines have a DE-loop of 8 residues, except for a small number of rabbit VH genes with DE-loop lengths 6 and 7, some of which are represented in the PDB (Table 2). All known mammalian VL germlines have a DE-loop of either 6 or 8 residues, except for one alpaca germline (IGLV5-12\*01) with a DE loop of length 3 (not represented in the PDB). There are 119 chains from 43 entries in the PDB with H4 loops longer than 8 amino acids, ranging from 10 amino acids to 16 amino acids. There are 65 chains from 50 entries that have insertions in L4 (all lambda chains), resulting in L4 loops of length 9. All of the antibodies in the PDB with non-germline insertions in L4 or H4, except one, are broadly neutralizing antibodies against HIV-1. The sole exception is an engineered nanobody against Hlg2 toxin (PDB: 5MJJE). Table 5 lists the various PDB structures of the bnAbs that have insertions in either the light or heavy chain DE loop, as well as their sequences, germlines, and their bnAb class. The structures are shown in Figure 8 and grouped by antibody class. In these structures, L4 and H4 bind the antigen epitope better than H1 and H3 as well as any of the light chain CDRs, as demonstrated by the extent of antigen-buried surface area for each CDR including CDR4 (Figure 9).

In the case of elongated H4 loops that contact gp120, the antibodies containing insertions in H4 bind to three separate binding sites. Antibodies in the VRC03/06/06B and 3BNC60/117 classes that have an Arg at position 1 of H4 that makes a salt bridge with residue Asp368 of gp120 near the CD4 binding site (Figure 8a), as noted in previous studies of the VRC01 class of antibodies.<sup>16,20</sup> Regardless of insertion length (4 residue insertion in the 3BNC antibodies, and 7 residue insertion in the VRC series of antibodies), these structures all localize to the same epitope on gp120 and share the same salt bridge as noted previously. In this way, DE residue 1 determines localization of the binding site of the VRC class of antibodies, which mimic the CD4 binding to gp120. Besides this interaction, much of the interaction between elongated H4 in the VRC antibodies and gp120 consists of buried hydrophobic contacts (Figure 8a). Figure 8b,c show the other structures of bnAbs with H4 loops with somatic insertions.

The non-canonical length L4 loops are related to the Hu\_IGLV3\_21\*01 germline and feature length 9 L4 loops.



**Figure 7.** Sequence entropy in naïve human antibodies and human germlines. **(a).** Sequence entropy for 12 common germlines in a naïve human antibody sample (>10,000 sequences for each germline). **(b).** Sequence entropy for human germlines derived from all IGKV, IGLV, and IGHV sequences from IMGT. From left to right, the pink shaded regions indicate CDR1, CDR2, and the DE loop. CDR3 is omitted due to varying lengths and different diversification mechanisms.

These loops not only directly bind antigen with hydrophobic interactions at the apex of the loop (Figure 8d), but also stabilize the conformation of L1 through a couple of backbone/backbone and backbone/side-chain hydrogen bonds to a serine in L1, which is sandwiched between L4 and L3. This ‘L1 sandwich’ motif appears to rigidify the binding conformation of the antibody light chain that buries a tremendous amount of binding surface area while binding to gp120 even in the presence of highly glycosylated elements (Figure 8d). The L1-14 conformations associated with these antibodies are

exclusive, and no other antibody structures contain these unique conformations of L1. The extended L4 loop length and conformation may stabilize the unusual conformation of L1, and enable the formation of new interactions between the antibody and antigen.

We aligned the DNA sequences of representatives of each class of antibodies in Table 5 with their germline sequences in order to identify the possible mechanism of the DE loop insertions. Somatic insertions are well studied in HIV bNAbs, and duplications of segments of the gene have been observed as

**Table 4.** Average sequence entropies for CDR and framework regions.

germline	CDR1	CDR2	FR1	FR2	FR3	CDR4
IGHV1-18*04	0.38	0.44	0.04	0.17	0.21	0.27
IGHV3-23*01	0.42	0.74	0.03	0.15	0.21	0.26
IGHV4-34*01	0.14	0.38	0.05	0.14	0.17	0.29
IGHV4-39*07	0.21	0.36	0.08	0.13	0.15	0.20
IGKV1-39*01	0.29	0.22	0.20	0.12	0.13	0.11
IGKV3-11*01	0.25	0.22	0.21	0.10	0.10	0.11
IGKV3-20*01	0.30	0.21	0.23	0.10	0.10	0.07
IGKV4-1*01	0.30	0.15	0.24	0.09	0.08	0.07
IGLV1-40*01	0.25	0.29	0.22	0.12	0.05	0.07
IGLV1-44*01	0.27	0.26	0.21	0.12	0.06	0.07
IGLV2-14*01	0.24	0.30	0.17	0.13	0.06	0.08
IGLV3-1*01	0.28	0.26	0.20	0.10	0.06	0.13
gene	CDR1	CDR2	FR1	FR2	FR3	CDR4
IGH	0.78	1.50	0.53	0.61	0.53	0.86
IGK	0.57	0.71	0.46	0.24	0.20	0.19
IGL	0.80	0.63	0.43	0.33	0.28	0.52

Average sequence entropies partitioned by CDR or framework region, excluding CDR3. Bolded values are those where the CDR4 average sequence entropy either compares to CDR1/CDR2, or exceeds the values for FR1, FR2, and FR3

the origin of some insertions.<sup>52</sup> These insertions are usually accompanied by somatic mutations either before or after the duplication event or both. Two of the antibodies listed in Table 5 appear to originate from such a duplication that maintains the reading frame. The antibody 35O22 contains a 24 base pair duplication of the DNA sequence that encodes the five amino acids before the H4 loop and the first three residues of the H4 loop (Figure 8e). The sequence identities of the duplicated segments to the germline are 21/24 base pairs (87.5%) and 15/24 base pairs (62.5%). Similarly, in the PGT122 gene, there is likely a 9 base pair duplication of the DNA segment that encodes the beta strand residues 4, 3, and 2 amino acids before the start of the DE loop (sequence FSG, Figure 8f). Structurally, the insertion alters the L4 loop and not the beta strand itself that is duplicated. The other 9 residue L4 loops all occur in the same germline (IGLV3-21\*01) and presumably occur via the same mechanism.

### Features of DE loop sequences from HIV-1 infected individuals

As noted earlier, DE loop insertions among 2.5 million antibody sequences from HIV-uninfected individuals are very rare. Of the entire data set from HIV-1 individuals, we found two unique heavy chain frame-shifted DE loops (IGHV1-18\*04; TNTPTSTA, IGHV5-10-1\*03; NKKSISTA), two unique kappa gene duplication insertions (IGKV3-20\*01; EKEKYRTDF, EKDKYRTD), and a single kappa gene frame-shifted DE loop (IGLV2-14\*02; SGGTAA). We also searched a set of ~24 million high-throughput sequences related to 13 studies of HIV-1 bnAbs to determine whether there were H4 and L4 insertions. These sequences are found across 13 HIV-1 high-throughput sequencing data sets related to the affinity maturation of VRC01, CH103, and PTG134-137 lineage antibodies, as well as co-evolution of HIV-1 bnAbs with their founder HIV-1 virus.<sup>17,46-51</sup>

We identified potential insertions in the DE loop by examining the alignments of these sequences to hidden Markov

models of the  $\kappa$ ,  $\lambda$ , and heavy chain variable domain, and searching for gaps in the HMM consensus sequence within 10 amino acids before and after the DE loop. This resulted in 599 unique (637 total) heavy chain sequences, 521 unique (1,354 total)  $\lambda$  sequences, and 3,174 (6,352 total)  $\kappa$  sequences with amino acid insertions in, or around the DE loop. We used Clustal-Omega to align each group of sequences (heavy chain,  $\kappa$ ,  $\lambda$ ) separately and Jalview to edit the alignments and to analyze the DE-loop region. In each of the alignments, there were some sequences that contained substantial changes in the protein sequence in or near the DE loop. These changes are likely due to frameshifts caused by an insertion or deletion of one or more bases that are compensated by another insertion or deletion later in the nucleic acid sequence such that the frame is restored. In some cases, these may be due to sequencing errors. In other cases, if the frameshift covers a substantial region of the amino acid sequence, the antibody may not fold properly.

Antibody L4 sequences with  $\kappa$  germlines (Figure 10a) have the most insertions compared to heavy and  $\lambda$  DE loop sequences. By comparing the sequence of 5 amino acids before the DE loop, the DE loop itself, and 5 amino acids after the DE loop, 480 of 3174 sequences (15.6%) appear to have a frameshift mutation in or adjacent to the DE loop. Of the remaining sequences, 98% contain a length 8 DE loop, with the sequence in 85% of these containing a two-residue insertion resulting in a sequence resembling GSGSGTDF (e.g., a sequence derived from IGHV4-4\*01 inserts GS before GSGTDF in Figure 10a). The remainder are of length 7. Some sequences have short insertions, alongside somatic mutations (e.g., a human IGKV3D-11\*01 sequence (germline sequence GPGTDF) has both a GS insertion before the DE loop, as well as eight-residue DE loop sequences ASAAGTEF and ASASGTDF).

Antibody L4 sequences from  $\lambda$  germlines (Figure 10b) have fewer insertions than  $\kappa$  L4 sequences, but still have features of simple amino-acid insertions, frameshifting insertions, and somatic mutation alongside insertion. Frameshifts are associated with 21.5% of the 521 unique sequences in the  $\lambda$  data set. Of the remaining 409 sequences, 92.4% have DE loops of length 8 and the rest are of length 7. Similar to IGHV4-4\*01, an IGLV3-12\*02 sequence has an inserted GS before the DE loop and mutations of the germline DE loop from NPGNTA, resulting in the sequence GSKSGNKA. Sequences from IGLV3-19\*01 and IGLV3-25\*03 have both frameshifts, as well as single amino-acid insertions that do not dramatically change the DE loop sequence (e.g., a single amino acid insertion changing the germline sequence from SSGNTA to STSGNTA).

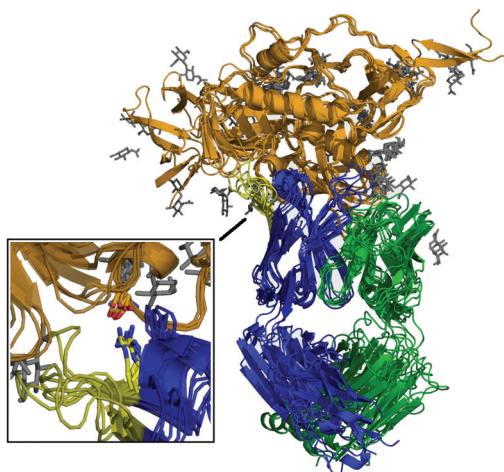
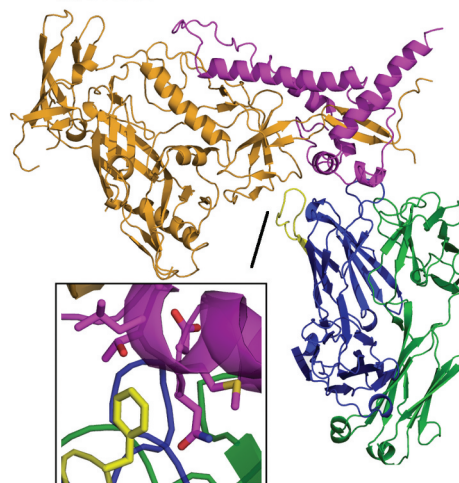
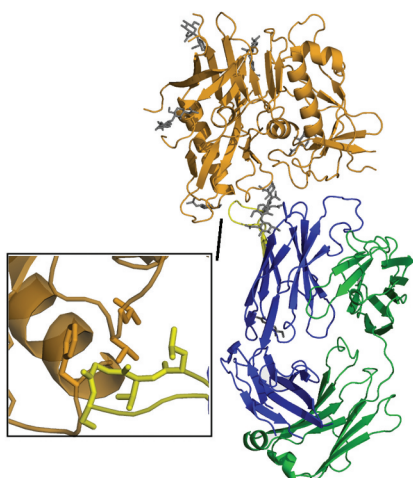
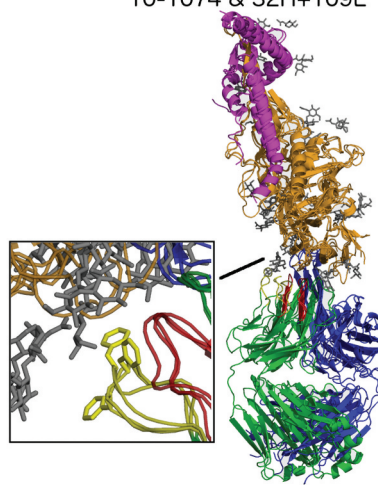
In the heavy-chain alignment, we observed many sequences that are likely due to frameshifts in or adjacent to the DE loop. 108 sequences (17.9%) have an insertion that causes a frameshift prior to the DE loop that extends through the DE loop sequence. 284 sequences (47.4%) contain a frameshift that starts prior to the DE loop and extends through the beginning of H3. The remaining 207 heavy chain sequences (34.5%) appear to have an insertion of one amino acid in the loop that leads into the d strand and no insertion in the DE loop. Examples of each case are shown in Figure 10c. There are

Table 5. HIV-1 bNAbs with insertions in L4 and/or H4.

Gene	PDB Chains	Length	DE loop sequence	Antibody/Germine
Heavy	4tyoYH 4tvpD* 5cezd* 5fyjD* 5fykD* 5fyjD* 5t3sD* 5u7oD* 5u7mD* 5um8D* 5utfd* 5utyD* 5v7jD* 5w6dD* 5wduH,M,U* 6ce0D* 6ch7D 6ch8D 6ch9D 6ck9D* 6de7D* 6ieqD 6mcoD 6mtjD 6mtdD 6mntnD 6mu6D 6mu7D 6mu8D 6mufD 6mugD 6nm6E 6nnfD 6nnjD 4jb9H	16	TDTEVPVTSFTSGAA RDT-----SISTA	35O22 IGHV1-18*01
Heavy	3se8H 4cc8F,H,I,5jxah 6cde5,Q,q* 6cdi8,Q,q* 6cue7,Q,q 6cuf8,Q,q 6espi,O,Y 6mpg8,Q,q 6mphQ,f,g 6n1vQ,f,g 6n1w8,Q,q 6nt2C,N, V 6osy8,F,P 6ot1J,S,q 6v8zC,I,O 4s1qh	15	RLFSQDLYPDRGTA RDT-----SISTA RQLSQDPDPDWGVA RDTST-----A	VRC06 IGHV1-2*02 VRC03 IGHV1-2*02
Heavy	4p9hH 4p9mH 5a7xN,P,R 5a8hF,L,R 5c7kE 5cjxA,D,H 5js9E 5jsaE 5thrP,R,T 5vykM,M,I 5vj6M,O,Q 6cm3P,R,T 6edup,R,T 6nqdC,G,K 4jpwH 4iswH 5v8G,H,I 5v8mH,I,R,S	15	RQLSQDPDPDWGIA RDTSIST-----A RQLSQDPDPDWGIA RDTSIST-----A RQLSQDPDPDWGIA RDTSIST-----A RQLSQDPDPDWGVA RDTSIST-----A	VRC01.H03 + 06.D-001739 IGHV1-2*02 N6 FR3-03 IGHV1-2*02 VRC01 FR303 IGHV1-2*02 VRC06B IGHV1-2*02
Heavy	6nm6U	15	AVDLTGSSPPIS	8ANC195
Heavy	6nnfU	15	ADESTST-----S	IGHV1-69*01
Heavy	4xnbZ,E,H	15	RHASWDFDYTSF	3BNC117
Heavy	4p9hH 4p9mH 5a7xN,P,R 5a8hF,L,R 5c7kE 5cjxA,D,H 5js9E 5jsaE 5thrP,R,T 5vykM,M,I 5vj6M,O,Q 6cm3P,R,T 6edup,R,T 6nqdC,G,K 4jpwH 4iswH 5v8G,H,I 5v8mH,I,R,S	12	RDTSIST-----A	IGHV1-2*02
Heavy	3rpiA,H 4gw4A,H	12	RQASWDFDYTSF	3BNC60
Lambda	4jy6A,C	9	RDTSIST-----A PDFRPGTTA	IGHV1-2*02 PGT123
Lambda	4fq2L 6ccbE,L 6ck9L* 5ceyA,C 5cezl* 5t3xL 5t3zL 5v7jL* 5w6dl* 6mtjL* 6mntL* 6mu6L* 6mu7L* 6mu8L* 6mufL* 6mugL* 6nm6L* 6nnfL* 6nnjL*	9	-----NSGNTA PDINFGTTA	IGLV3-21*01 10-1074
Lambda	4r26L 4r2gC,I,M 5t3sL* 5um8L* 6ce0L* 6ieqL* 6mcoL* 6mdtL*	9	-----NSGNTA	IGLV3-21*01
Lambda	5cexC	9	PDINFGTTA -----NSGNTA	PGT124 IGLV3-21*01
Lambda	4fq1L 4fqcl 4jy4A	9	PDSNFGTTA -----NSGNTA	32H + 109L IGLV3-21*01
Lambda	3jcb8 3jcc8 4jy5L 4ncoG,K,C 4tvpL* 5d9qL,E,M 5fyjL* 5fykL* 5fyjL* 5i8hJ,L 5u7mL* 5u7oL* 5utfl* 5utyL* 5wduB,K,S* 6b0nL 6cdeN,8,n* 6cdiN,6,n* 6cuf6,n,n* 6de7L* 6nt2F,T,K* 6osyD,N,6* 6ot1R,I,n*	9	PDSPFGTTA -----NSGNTA PGSTFGTTA -----NSGNTA	PGT121 IGLV3-21*01 PGT122 IGLV3-21*01

The table lists all of the antibody structures in the PDB with insertions in L4 or H4 related to bNAbs along with their sequences, germline, and bNA lineage. Entries with an asterisk (\*) represent structures with insertions in both the light and heavy chains. There is only one other antibody with an inserted DE loop in the PDB: the engineered nanobody toward HlgB2 toxin in cholera (PDBID 5mje, DE sequence RDSAEDSAKNTV). It is not listed in the Table. The amino acid sequences of the inserted DE loop in the PDB and the germline were aligned by locating the insertion in their DNA sequences.



**a.** VRC-03/06/06B & 3BNC-60/117**b.** 35022**c.** 8ANC195**d.** PTG-121/2/3/4 & 10-1074 & 32H+109L

**e.** DE loop

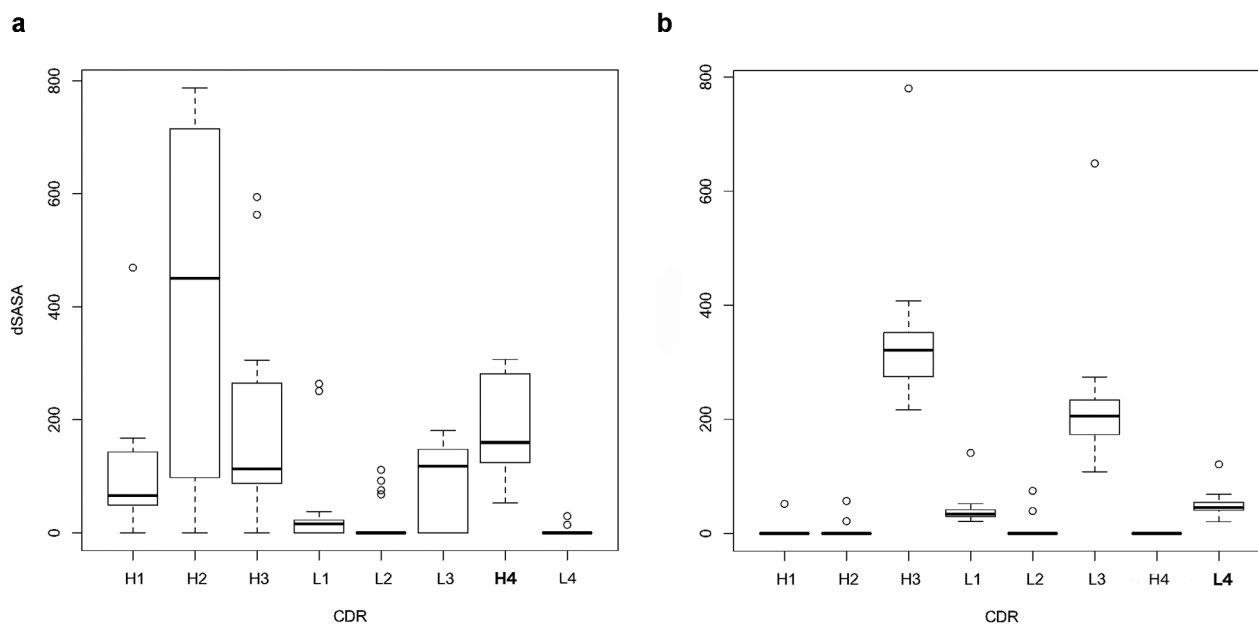
IGHV1-18*02	K	L	Q	G	<u>R</u>	<u>V</u>	<u>T</u>	<u>M</u>	<u>T</u>	<u>T</u>	<u>D</u>	<u>T</u>				<u>S</u>	<u>T</u>	<u>S</u>	<u>T</u>	<u>A</u>					
IGHV1-18*02	aag	ctc	cag	ggc	<u>aga</u>	<u>gtc</u>	<u>acc</u>	<u>atg</u>	<u>acc</u>	<u>aca</u>	<u>gac</u>	<u>aca</u>	<u>aga</u>	<u>gtc</u>	<u>acc</u>	<u>atg</u>	<u>acc</u>	<u>aca</u>	<u>gac</u>	<u>aca</u>	tcc	acg	agc	aca	gcc
35022	gcc	ttt	caa	gac	<u>aga</u>	<u>gtc</u>	<u>att</u>	<u>atg</u>	<u>acg</u>	<u>aca</u>	<u>gac</u>	<u>aca</u>	<u>gaa</u>	<u>gtc</u>	<u>cct</u>	<u>gtg</u>	<u>acc</u>	<u>tca</u>	<u>ttc</u>	<u>acg</u>	tcc	acg	ggc	gca	gcc
35022	A	F	Q	D	<u>R</u>	<u>V</u>	<u>I</u>	<u>M</u>	<u>T</u>	<u>T</u>	<u>D</u>	<u>T</u>	<u>E</u>	<u>V</u>	<u>P</u>	<u>V</u>	<u>T</u>	<u>S</u>	<u>F</u>	<u>T</u>	<u>S</u>	<u>T</u>	<u>G</u>	<u>A</u>	<u>A</u>
DE loop					<u>T</u>	<u>D</u>	<u>T</u>	<u>E</u>	<u>V</u>	<u>P</u>	<u>V</u>	<u>T</u>	<u>S</u>	<u>F</u>	<u>T</u>	<u>S</u>	<u>T</u>	<u>G</u>	<u>A</u>	<u>A</u>					

**f.** DE loop

IGLV3-21*01	I	P	E	R	<u>F</u>	<u>S</u>	<u>G</u>	<u>F</u>	<u>S</u>	<u>G</u>	<u>S</u>	<u>N</u>	<u>S</u>	<u>G</u>	<u>N</u>	<u>T</u>	<u>A</u>								
IGLV3-21*01	atc	cct	gag	cga	<u>ttc</u>	<u>tct</u>	<u>ggc</u>	<u>ttc</u>	<u>tct</u>	<u>ggc</u>	tcc	aac	tct	ggg	aac	acg	gcc								
PGT122	att	cct	gac	cga	<u>ttt</u>	<u>tct</u>	<u>ggg</u>	<u>tcc</u>	<u>cct</u>	<u>ggc</u>	tcc	act	ttt	ggg	acc	acg	gcc								
PGT122	I	P	D	R	<u>F</u>	<u>S</u>	<u>G</u>	<u>S</u>	<u>P</u>	<u>G</u>	<u>S</u>	<u>T</u>	<u>F</u>	<u>G</u>	<u>T</u>	<u>T</u>	<u>A</u>								
DE loop					<u>P</u>	<u>G</u>	<u>S</u>	<u>T</u>	<u>F</u>	<u>G</u>	<u>T</u>	<u>T</u>	<u>A</u>												

**Figure 8.** Alignment of a subset of gp120 binding HIV-1 bNAb s representing all unique DE loop sequences. (a). Aligned structures of H4 inserted bNAb s of the VRC and 3BNC series (one structure per H4 sequence). The inset shows hydrophobic contacts between the antibody-antigen interface, as well as a salt bridge between the first Arg residue in H4 and the antigen. (b). Structure of the 35022 antibody found in 35 PDB entries with an H4 of length 16. The inset shows a hydrophobic residue interaction with hydrophobic residues in the helix of the antigen. (c). Structure of the 8ANC195 antibody with an H4 of length 12 represented in 16 PDB entries (one structure per H4 sequence). The inset shows hydrophobic residues interacting with hydrophobic residues in a loop of the antigen. (d). Aligned structures of L4-inserted bNAb s binding to HIV-1 gp120 (one representative per unique L4 sequence of length 9). The inset shows hydrophobic contacts at the antibody-antigen interface, including hydrophobic contacts to glycosylation residues on gp120 (gray sticks). (e). Alignment of antibody 35022 heavy chain DNA and protein sequences around the DE loop segment with the DNA and protein sequences of the germline IGHV1-18\*02. The 35022 DE loop protein sequence is at the bottom in yellow. The IGHV1-18\*02 DE loop protein sequence is at top. The duplicated regions are underlined. The IGHV1-18\*02 is identical in both duplicated regions, while the 35022 sequence has diverged either prior to or after duplication or both. Identical base pairs or amino acids between 35022 and IGHV1-18\*02 are highlighted in gray. (f). Alignment of antibody PGT122 light chain DNA and protein sequences around the DE loop segment with the DNA and protein sequences of the germline IGLV3-21\*01. The PGT122 DE loop protein sequence is at the bottom in yellow. The IGLV3-21\*01 DE loop protein sequence is at top. The duplicated regions are underlined. The IGLV3-21\*01 is identical in both duplicated regions, while the PGT122 sequence has diverged either prior to or after duplication or both. Identical base pairs or amino acids between PGT122 and IGLV3-21\*01 are highlighted in gray.





**Figure 9.** Buried surface area for each CDR at the antibody-antigen interface of HIV-1 bNAbs that bind to gp120 in the PDB. (a). Buried surface area plot for 18 PDB structures (non-redundant by chain) with insertions in H4. (b). Buried surface area plot for 31 PDB structures (non-redundant by chain) with insertions in L4.

no clear examples of a simple insertion within the DE loop itself in the heavy chain sequences, and no cases that resemble the insertions in the HIV bNAbs in the PDB.

## Discussion

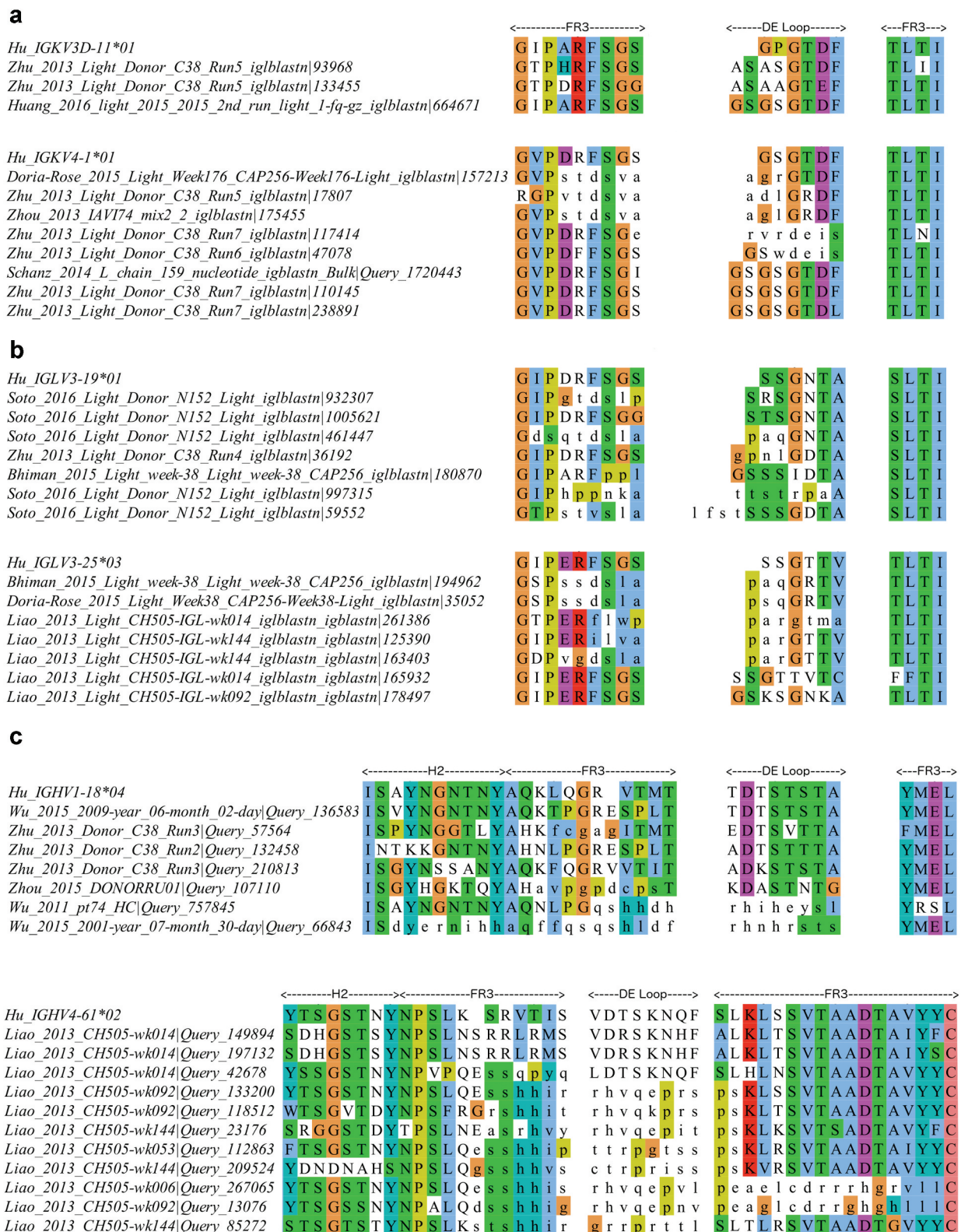
In this study, we analyzed sequence and structural features of the antibody DE loop. We clustered DE loop conformations of the heavy and light chains, identified atomic interactions that are highly associated with various CDR1/CDR2 and DE loop pairs, and have shown features of affinity matured antibodies that use DE loops with somatic insertions to directly bind antigen. We identified nine clusters of the DE loop in heavy and light chains, and denoted them L4-6-1, L4-6-2, L4-6-3, L4-6-4, L4-8-1, H4-6-1, H4-7-1, H4-8-1, and H4-8-2. For each clustering of each length of L4 and H4, we also provided support for each individual cluster by counting how many structures and unique sequences after enforcing a strict electron density fit calculation through the EDIA score.<sup>15</sup> This allowed us to find clusters that are not defined by mis-solved residues, such as peptide flips,<sup>55</sup> which is an important step in validating clusters and backbone structures not only for clustering the DE loop, but the other CDRs as well as backbone structures from other proteins.<sup>64,65</sup>

With this new structural classification and nomenclature to describe the DE loop of antibody structures, we argue that the DE loop structure and sequence should be analyzed when new antibody structures are determined and during development of antibody reagents and therapeutics. In some respects, the DE loop acts as a fourth CDR, since it behaves like the other CDRs: it binds antigen and it undergoes somatic mutations and insertions that directly or indirectly affect the antibody binding paratope at a rate that is higher than most framework residues (at least in the heavy chain and  $\lambda$  light-chain sequences).

Our analysis has the limitation that we cannot always determine causative associations between the sequence and structure of the DE loop and those of the CDRs. For example, even if we have some association between a less common DE loop conformation and CDR1 cluster, we cannot say that DE loop conformation is a determinant of the CDR1/CDR2 conformation. For example, L1-13-1 conformations are associated with L4-6-2, while L1-13-2 conformations are associated with L4-8-1. The sequence profiles of L1-13-1 and L1-13-2 are distinct. We do not know whether the L1-13-2 conformation would be maintained if an antibody was constructed with an L1-13-2 sequence and a typical L4-6-2 DE loop. Therefore, when accounting for effects that the DE loop may have on a neighboring CDR, it is important to note the L4 conformations available for that particular choice of L1 conformation, and also the sequence positions within a DE loop cluster that differentially affect CDR1 or CDR2 conformation. While some hydrogen bonds occur at a very high occupancy (90% or above) in specific DE loops pairings with CDR1/CDR2, others occur at a lower occupancy. Hydrogen bonds of this nature may not be a strict requirement for a particular L1 conformation, but may have an effect on stability of that L1 conformation, but, without extensive additional experimental data to test the hypothesis, we cannot determine this.

Considering the DE loop as a fourth CDR suggests applications for antibody design and antibody modeling. For example, when designing antibodies using the ‘CDR grafting’ method,<sup>66</sup> we suggest that whenever CDR1 is grafted on the light chain, or CDR1 or CDR2 on the heavy chain, L4 or H4 should be ‘co-grafted’ onto the same template structure. This method will preserve contacts between L4/L1, H4/H1, or H4/H2 that may be necessary for preserving or stabilizing the structures of CDR1 or CDR2.

When considering antibody structure prediction,<sup>67–69</sup> a common strategy is to use CDR and framework templates



**Figure 10.** DE loop and DE loop adjacent insertions from a large antibody sequencing data set from HIV-infected individuals. (a). Insertions in  $\kappa$  gene antibodies. (b). Insertions in  $\lambda$  gene antibodies. C. Insertions in heavy gene antibodies.

based upon sequence similarity to known structures. We suggest extra attention to the relationships of L4 sequences with their structural clusters. For example,  $\kappa$  antibodies with a somatic mutation at the first position of the L4 from glycine

to any other residue should be modeled with representative structures from cluster L4-6-2 instead of the more common L4-6-1 conformation. A similar approach can be considered when selecting templates for solving X-ray crystallographic

structures by molecular replacement. Taking this information into account is more likely to recapitulate contacts observed in experimental structures. The appropriate cluster, and thus structure, for CDR1 and CDR2 often depends on the sequence and conformation of CDR4, and they should be modeled together in antibody structure prediction methods.

With high-throughput sequencing data in response to HIV-1, we have shown that the DE loop undergoes somatic mutation, alongside nucleotide insertions and deletions causing frameshift mutations in several human germline examples. Tracking useful features from the DE loop sequences that contribute to antigen binding, and ultimately neutralization of viral infections, may prove an important step in identifying functional antibodies from the human repertoire.

## Materials and methods

### Antibody structure and sequence data

We compiled sequence and structure data for all antibodies from the PDB. To collect the list of antibodies in the PDB, we used a lab-maintained software, PyIgClassify.<sup>70</sup> PyIgClassify compiles all antibody structures from the PDB by applying a set of hidden Markov models (HMMs) for each antibody gene (VH, V $\lambda$ , and V $\kappa$ ) to all sequences in the PDB using HMMER3.0. PyIgClassify also rennumbers antibodies according to a modified Honegger-Plückthun CDR scheme and numbering system described in North et al.<sup>56,71</sup> In order to identify CDRs in PyIgClassify, the software uses sequence alignment to the match states of the HMMs.

In order to identify which residues are structurally variable, we plotted  $\phi$  and  $\psi$  for all residues in and around the solvent-exposed DE loop (3 before the loop, and 3 after the loop, Figure 2). We updated PyIgClassify to recognize L4 and H4 in each antibody sequence, adding insert codes appropriately for loops longer the pre-allocated range of numbers.

We determine germline by comparing each PDB sequence to a curated set of IMGT germline protein sequences with BLAST, taking into account the author-provided species designation. However, these are often incorrect. We use a germline from a different species from the author-provided one if the sequence identity of the antibody is at least 8 percentage points higher than the author-provided species. This script also handles cases of ambiguous assignment, such as humanized antibodies originating from non-human germlines. The data set used in this study was compiled in March 2020, and includes data for 3,910 antibody structures containing 13,012 individual chains (available at <http://dunbrack2.fccc.edu/PyIgClassify>)

### Calculating electron density support for individual atoms

In order to add support for the electron density for the individual atoms for all structures within clusters, we calculated the EDIA score using the *ediascorer* standalone software from the University of Hamburg (<https://www.zbh.uni-hamburg.de/en/forschung/amd/software/edia.html>).<sup>15</sup> This software requires the structure factor file as well as the electron density map (.mtz). We downloaded both of these files from the PDB except

for some older entries without deposited structure factor files. We then calculated the EDIA score for each backbone atom, and took as the EDIA score the value for the carbonyl oxygen, which is the most sensitive atom to the electron density fit. We established a second data set by eliminating any structure with one or more carbonyl oxygen atoms with an EDIA score of less than 0.75.

### Analyzing antibody-antigen complex set

For non-canonical length structures, we calculated the antibody-antigen buried surface area with the Rosetta macromolecular modeling suite.<sup>72</sup> We calculated buried surface area as the change in antigen surface area of the CDR from the bound structure to unbound structure:

$$dSASA = SA_{bound} - SA_{unbound} \quad (\text{eq.1})$$

where SA represents the surface area calculated in Rosetta using the Shrake-Rupley algorithm<sup>73</sup> and a standard probe radius of 1.4 Å.

### Clustering loop structures

In order to group various conformations of L4 and H4 into structural families, we implemented a density-based clustering method for dihedral angles based on the DBSCAN algorithm.<sup>53</sup> This unsupervised learning method finds robust clusters by identifying dense regions in the metric space which are separated by low density. It also identifies “noise points,” which are outlier structures due to poor crystal structure determination or unusual mutations that cause uncommon structural changes. We used the implementation of DBSCAN in the scikit learn library in Python.<sup>74</sup>

To compare two loops  $i$  and  $j$  with identical lengths, we first calculate the dihedral similarity between two angles  $\theta_1$  and  $\theta_2$  for each pair of corresponding residues, where  $\theta$  represents any chosen backbone dihedral angle selected from  $\phi$ ,  $\psi$ , or  $\omega$ :

$$d = 2(1 - \cos \theta_1 - \theta_2) \quad (\text{eq.2})$$

For our purposes we chose to include  $\phi$ ,  $\psi$ , and  $\omega$ , which provides the maximum capability to resolve structures with both cis- and trans- peptide bonds. Next, we take as the final clustering distance the maximum value out of the set of calculations of  $d$  for  $\{\phi, \psi, \omega\}$ , which we call the  $L_\infty$  norm:

$$L_\infty = \text{Max}(d_{\phi_1}, d_{\psi_1}, d_{\omega_1}, \dots, d_{\phi_N}, d_{\psi_N}, d_{\omega_N}) \quad (\text{eq.3})$$

We chose the  $L_\infty$  norm due to its sensitivity in separating loops which are different even at one single dihedral, giving our final clustering single dihedral resolution.

The resulting set of pairwise  $L_\infty$  distances are then clustered from an  $N \times N$  pairwise matrix using DBSCAN. This algorithm requires two parameters:  $\epsilon$  and *MinPts*. The first parameter,  $\epsilon$ , describes a distance from a given data point to search for neighboring data points. The second parameter, *MinPts*, specifies the requirement for the minimum number of neighboring data points within  $\epsilon$  of a data point to label the data point under consideration a ‘core point’. Data points that are within  $\epsilon$  of a core point, but do not have *MinPts* data points within  $\epsilon$  are called



'border points'; points that do not meet either criterion are labeled as 'noise points.' The final clusters are the connected graphs of all of the core points, together with their border points.

Each selection of a combination of  $\epsilon$  and *MinPts* produces a different set of clusters. Two main obstacles exist in identifying all of the interesting clusters from DBSCAN. First, at larger values of  $\epsilon$  and smaller *MinPts*, DBSCAN may merge clusters that ought to be separated. Merged clusters are easily identified from their Ramachandran plot distributions at specific residues within a cluster (e.g., separated densities in the alpha and beta regions of the Ramachandran map). Second, valid, low-density clusters may only arise at larger values  $\epsilon$ , while also producing undesirable merged clusters. This means that no singular selection of  $\epsilon$  and *MinPts* will generate the entire set of valid clusters. To overcome these two issues, we developed a method to select a set of final clusters after running DBSCAN on a grid of  $\epsilon$  and *MinPts*, by combining the results of each run of DBSCAN. First, we establish a parameter grid of  $\epsilon$  and *MinPts* by selecting a range of both parameters, and run DBSCAN at each parameter selection. We then filter out any merged clusters by removing any clusters in which any two members of the cluster are more than 150° apart. This eliminates clusters that contain points in two or more Ramachandran regions (A, B, L, or E), the centroids of which are approximately 180° apart in  $\phi$  or  $\psi$ . Next, the remaining clusters that pass the merge filtering criterion are treated as nodes on a graph, where the nodes have edges connected to them based on the calculation of Simpson's similarity score:<sup>75</sup>

$$S = \frac{|A \cap B|}{\min(|A|, |B|)} \quad (\text{eq.4})$$

Finally, for each connected subgraph with  $n$  nodes, we take the final cluster of that subgraph as the union of all nodes  $n$  within the connected subgraph. This produces a final clustering set with clusters of varying density, without including merged clusters.

Following the determination of the final cluster set, we determined cluster representatives using angular statistical analysis. For a given cluster C consisting of N data points, for each structure  $i$  we calculate the average distance  $d_i$  to all other points  $j$  in the same cluster C:

$$d_i = \frac{1}{N} \sum_{j=1, N} d_{ij} \quad (\text{eq.5})$$

We choose the cluster representative as the structure that has the lowest  $d_i$  of all of the structures.

### Identifying important hydrogen bonds between CDR4 and CDR1/CDR2

We calculated all hydrogen bonds between CDR4, CDR1, and CDR2 using Rosetta's distance and orientation-dependent hydrogen bond energy accessed through the *report\_hbonds\_for\_plugin.<release>* available in the public release of Rosetta3. We used the resulting contact information to find important contacts that are either frequent or unique over several CDR-lengths and germlines. We analyzed the hydrogen bonds between all CDR1-CDR4 and CDR2-CDR4 pairs for which

both CDR1 and CDR4 have defined cluster membership. We then calculated the hydrogen bond occupancy for a particular hydrogen bond as the following:

$$\text{occupancy} = \frac{\#\text{hbondstocdr}|\text{deresidue, deatom, cdrcluster}}{\#\text{structures}|\text{deresidue, deatom, cdrcluster}} \quad (\text{eq.6})$$

### High-throughput sequence analysis of naïve human antibodies

We accessed high-throughput sequencing data through the antibodymap.org server ([www.antibodymap.org](http://www.antibodymap.org)).<sup>76</sup> To gain an understanding of how variable L4 and H4 are compared to the other CDRs, we analyzed 12 human germlines (IGHV1-18\*04, IGHV3-23\*01, IGHV4-34\*01, IGHV4-39\*07, IGKV1-39\*01, IGKV3-11\*01, IGKV3-20\*01, IGKV4-1\*01, IGLV1-40\*01, IGLV1-44\*01, IGLV2-14\*01, IGLV3-1\*01) collected from naïve donor deep sequencing samples with thousands of sequences for each germline (download shell script included in supplementary data). Separately, we compared sequence variability between all human germlines for each heavy,  $\lambda$ , and  $\kappa$  gene compiled from IMGT for all germline sequences of the same length. We calculated sequence variability according to the Shannon entropy, denoted  $H$ , which represents a robust method to calculate antibody CDR variability:<sup>77</sup>

$$H = - \sum_{i=1, N} p_i \log_2 p_i \quad (\text{eq.7})$$

We calculated  $H$  only for residues up until the conserved cysteine before CDR3 on both the light and heavy chains.

### High-throughput sequence analysis of HIV-1 bnAbs

In order to search for insertions in L4 or H4 amongst HIV-1 infected patients, we collected all studies referring to HIV-1 from the antibodymap API (download shell script included in supplementary data). We identified CDRs for all of the FASTA files using the HMMER3.0 *hmmsearch* command, providing the profile HMMs implemented in PyIgClassify for IGHV, IGKV, IGLV, and IGLV6 genes (provided in supplementary data). We searched for sequences that had insertions compared to the profile, and examined these for features related to the long L4 or H4 structures we found in the PDB (sequences are provided in FASTA format in the supplementary data). From this set, we used Clustal-omega to align all of the sequences to the germline sequence which matched the IMGT germline assignment (provided in supplementary data). We observed frameshift mutations using the IMGT/V-QUEST tool,<sup>78</sup> which notates nucleotide insertions that result in frameshifts.

### Funding

This work was funded by NIH grant [R35 GM122517 (to R.L.D.)], NIH grant [P30 CA006927 (to Fox Chase Cancer Center)], and NIH Structural Biology and Molecular Biophysics training grant awarded by the Department of Biochemistry and Molecular Biophysics at the University of Pennsylvania [T32 GM008275 (to S.P.K.)].

## References

- Kabat EA, Wu TT. Attempts to locate complementarity-determining residues in the variable positions of light and heavy chains \*. *Ann N Y Acad Sci.* 1971;190:382–93. doi:10.1111/j.1749-6632.1971.tb13550.x.
- Bork P, Holm L, Sander C. The immunoglobulin fold. *J Mol Biol.* 1994;242:309–20. doi:10.1006/jmbi.1994.1582.
- Lesk AM, Chothia C. Evolution of proteins formed by  $\beta$ -sheets: II. The core of the immunoglobulin domains. *J Mol Biol.* 1982;160:325–42. doi:10.1016/0022-2836(82)90179-6.
- Lefranc M-P, Giudicelli V, Ginestoux C, Bodmer J, Müller W, Bontrop R, Lemaitre M, Malik A, Barbié V, Chaume D. IMGT, the international ImMunoGeneTics database. *Nucleic Acids Res.* 1999;27:209–12. doi:10.1093/nar/27.1.209.
- Chothia C, Lesk AM. Canonical structures for the hypervariable regions of immunoglobulins. *J Mol Biol.* 1987;196:901–17. doi:10.1016/0022-2836(87)90412-8.
- Tramontano A, Chothia C, Lesk AM. Framework residue 71 is a major determinant of the position and conformation of the second hypervariable region in the VH domains of immunoglobulins. *J Mol Biol.* 1990;215:175–82. doi:10.1016/S0022-2836(05)80102-0.
- Foot J, Winter G. Antibody framework residues affecting the conformation of the hypervariable loops. *J Mol Biol.* 1992;224:487–99. doi:10.1016/0022-2836(92)91010-M.
- Al-Lazikani B, Lesk AM, Chothia C. Standard conformations for the canonical structures of immunoglobulins 11 Edited by I. A. Wilson. *J Mol Biol.* 1997;273:927–48. doi:10.1006/jmbi.1997.1354.
- Carter P, Presta LEN, Gorman CM, Ridgway JB, Henner D, Wong WL, Rowland AM, Kotts C, Carver ME, Shepard HM. Humanization of an anti-p185HER2 antibody for human cancer therapy. *Proc Natl Acad Sci.* 1992;89:4285–89.
- Hsiao K, Bajorath J, Harris LJ. Humanization of 60.3, an anti-CD18 antibody; importance of the L2 loop. *Protein Eng Des Sel.* 1994;7:815–22. doi:10.1093/protein/7.6.815.
- Krauss J, Arndt MAE, Zhu Z, Newton DL, BK V, Choudhry V, Darbha R, Ji X, NS C-L, MP D, et al. Impact of antibody framework residue V H –71 on the stability of a humanised anti-MUC1 scFv and derived immunoenzyme. *Br J Cancer.* 2004;90:1863–70. doi:10.1038/sj.bjc.6601759.
- Maruthachalam BV, El-Sayed A, Liu J, AR S, Hill W, MK A, Pastushok L, Fonge H, Barreto K, Geyer CR. A single-framework synthetic antibody library containing a combination of canonical and variable complementarity-determining regions. *ChemBioChem.* 2017;18:2247–59. doi:10.1002/cbic.201700279.
- Teplyakov A, Obmolova G, Malia TJ, Raghunathan G, Martinez C, Fransson J, Edwards W, Connor J, Husovsky M, Beck H, et al. Structural insights into humanization of anti-tissue factor antibody 10H10. *mAbs.* 2018;10:269–77. doi:10.1080/19420862.2017.1412026.
- Lehmann A, Wixted Josephine H F, Shapovalov MV, Roder H, Dunbrack RL, Robinson MK. Stability engineering of anti-EGFR scFv antibodies by rational design of a lambda-to-kappa swap of the VL framework using a structure-guided approach. *mAbs.* 2015;1–14. doi:doi:doi:10.1080/19420862.2015.1088618 6 7 45–1071
- Meyder A, Nittinger E, Lange G, Klein R, Rarey M. Estimating electron density support for individual atoms and molecular fragments in X-ray structures. *J Chem Inf Model.* 2017;57:2437–47. doi:10.1021/acs.jcim.7b00391.
- Wu X, Zhou T, Zhu J, Zhang B, Georgiev I, Wang C, Chen X, Longo NS, Louder M, McKee K, et al. Focused evolution of HIV-1 neutralizing antibodies revealed by structures and deep sequencing. *Science.* 2011;333:1593–602. doi:10.1126/science.1207532.
- Zhou T, Zhu J, Wu X, Moquin S, Zhang B, Acharya P, Georgiev IS, Altae-Tran HR, Chuang G-Y, Joyce MG, et al. Multidonor analysis reveals structural elements, genetic determinants, and maturation pathway for HIV-1 neutralization by VRC01-class antibodies. *Immunity.* 2013;39:245–58. doi:10.1016/j.immuni.2013.04.012.
- Huang J, Kang BH, Ishida E, Zhou T, Griesman T, Sheng Z, Wu F, Doria-Rose NA, Zhang B, McKee K, et al. Identification of a CD4-Binding-Site Antibody to HIV that Evolved Near-Pan Neutralization Breadth. *Immunity.* 2016;45:1108–21. doi:10.1016/j.immuni.2016.10.027.
- Rutten L, Lai Y-T, Blokland S, Truan D, Bisschop IJM, Strokappe NM, Koornneef A, Manen DV, Chuang GY, Farney SK, et al. A universal approach to optimize the folding and stability of prefusion-closed HIV-1 envelope trimers. *Cell Rep.* 2018;23:584–95. doi:10.1016/j.celrep.2018.03.061.
- Scheid JF, Mouquet H, Ueberheide B, Diskin R, Klein F, Oliveira TYK, Pietzsch J, Fenyo D, Abadir A, Velinon K, et al. Sequence and structural convergence of broad and potent HIV antibodies that mimic CD4 binding. *Science.* 2011;333:1633–37.
- Mouquet H, Scharf L, Euler Z, Liu Y, Eden C, Scheid JF, Halper-Stromberg A, Gnanapragasam PNP, Spencer DIR, Seaman MS, et al. Complex-type N-glycan recognition by potent broadly neutralizing HIV antibodies. *Proc Natl Acad Sci.* 2012;109:E3268–E3277. doi:10.1073/pnas.1217207109.
- Klein F, Diskin R, Scheid JF, Gaebler C, Mouquet H, Georgiev IS, Pancera M, Zhou T, Icesu R-B, Fu BZ, et al. Somatic mutations of the immunoglobulin framework are generally required for broad and potent HIV-1 neutralization. *Cell.* 2013;153:126–38. DOI:10.1016/j.cell.2013.03.018.
- Julien JP, Sok D, Khayat R, Lee JH, Doores KJ, Walker LM, Ramos A, Diwanji DC, Pejchal R, Cupo A, et al. Broadly neutralizing antibody PGT121 allosterically modulates CD4 binding via recognition of the HIV-1 gp120 V3 base and multiple surrounding glycans. *PLoS Pathog.* 2013;9:e1003342–e1003342. doi:10.2210/pdb4jy4/pdb.
- Julien J-P, Cupo A, Sok D, Stanfield RL, Lyumkis D, Deller MC, Klasse P-J, Burton DR, Sanders RW, Moore JP, et al. Crystal Structure of a Soluble Cleaved HIV-1 Envelope Trimer. *Science.* 2013;342:1477–83. doi:10.1126/science.1245625.
- Scharf L, Scheid JF, Lee JH, West AP, Chen C, Gao H, Gnanapragasam PNP, Mares R, Seaman MS, Ward AB, et al. Antibody 8ANC195 reveals a site of broad vulnerability on the HIV-1 envelope spike. *Cell Rep.* 2014;7:785–95. doi:10.1016/j.celrep.2014.04.001.
- Garces F, Sok D, Kong L, McBride R, Kim HJ, Saye-Francisco KF, Julien J-P, Hua Y, Cupo A, Moore JP, et al. Structural evolution of glycan recognition by a family of potent HIV antibodies. *Cell.* 2014;159:69–79. doi:10.1016/j.cell.2014.09.009.
- Pancera M, Zhou T, Druz A, Georgiev IS, Soto C, Gorman J, Huang J, Acharya P, Chuang G-Y, Ofek G, et al. Structure and immune recognition of trimeric pre-fusion HIV-1 Env. *Nature.* 2014;514:455–61. doi:10.1038/nature13808.
- Wu X, Zhang Z, Schramm CA, Joyce MG, Do Kwon Y, Zhou T, Sheng Z, Zhang B, O'Dell S, McKee K, et al. Maturation and diversity of the VRC01-antibody lineage over 15 years of chronic HIV-1 infection. *Cell.* 2015;161:470–85. doi:10.1016/j.cell.2015.03.004.
- Scharf L, Wang H, Gao H, Chen S, McDowall AW, Bjorkman PJ. Broadly neutralizing antibody 8ANC195 recognizes closed and open states of HIV-1. *Env Cell.* 2015;162:1379–90. doi:10.1016/j.cell.2015.08.035.
- Kong L, Torrents de la Peña A, Deller MC, Garces F, Sliopen K, Hua Y, Stanfield RL, Sanders RW, Wilson IA. Complete epitopes for vaccine design derived from a crystal structure of the broadly neutralizing antibodies PGT128 and 8ANC195 in complex with an HIV-1 Env trimer. *Acta Crystallogr D Biol Crystallogr.* 2015;71:2099–108. doi:10.1107/S1399004715013917.
- Kong R, Xu K, Zhou T, Acharya P, Lemmin T, Liu K, Ozorowski G, Soto C, Taft JD, Bailer RT, et al. Fusion peptide of HIV-1 as a site of vulnerability to neutralizing antibody. *Science.* 2016;352:828–33. doi:10.1126/science.aae0474.
- Davenport TM, Gorman J, Joyce MG, Zhou T, Soto C, Guttman M, Moquin S, Yang Y, Zhang B, Doria-Rose NA, et al. Somatic hypermutation-induced changes in the structure and dynamics of HIV-1 broadly neutralizing antibodies. *Structure.* 2016;24:1346–57. doi:10.1016/j.str.2016.06.012.



33. Steichen JM, Kulp DW, Tokatlian T, Escolano A, Dosenovic P, Stanfield RL, McCoy LE, Ozorowski G, Hu X, Kalyuzhnyi O, et al. HIV vaccine design to target germline precursors of glycan-dependent broadly neutralizing antibodies. *Immunity*. 2016;45:483–96. doi:10.1016/j.immuni.2016.08.016.
34. Wang H, Cohen AA, Galimidi RP, Gristick HB, Jensen GJ, Bjorkman PJ. Cryo-EM structure of a CD4-bound open HIV-1 envelope trimer reveals structural rearrangements of the gp120 V1V2 loop. *Proc Natl Acad Sci*. 2016;113:E7151–E7158. doi:10.1073/pnas.1615939113.
35. Pancera M, Lai Y-T, Bylund T, Druz A, Narpala S, O'Dell S, Schön A, Bailer RT, Chuang G-Y, Geng H, et al. Crystal structures of trimeric HIV envelope with entry inhibitors BMS-378806 and BMS-626529. *Nat Chem Biol*. 2017;13:1115–22. doi:10.1038/nchembio.2460.
36. Guenaga J, Garces F, Stanfield RL, Dubrovskaya V, Higgins B, Carrette B, Ward AB, Wilson IA, Wyatt RT. Glycine substitution at helix-to-coil transitions facilitates the structural determination of a stabilized subtype C HIV envelope glycoprotein. *Immunity*. 2017;46:792–803.e3. doi:10.1016/j.immuni.2017.04.014.
37. Chuang G-Y, Geng H, Pancera M, Xu K, Cheng C, Acharya P, Chambers M, Druz A, Tsybovsky Y, Wanninger TG, et al. Structure-based design of a soluble prefusion-closed HIV-1 env trimer with reduced CD4 affinity and improved immunogenicity. *J Virol*. 2017;91:e02268–16. doi:10.1128/JVI.02268-16.
38. Zhou T, Doria-Rose NA, Cheng C, GBE S-J, Chuang G-Y, Chambers M, Druz A, Geng H, McKee K, Kwon YD, et al. Quantification of the impact of the HIV-1-glycan shield on antibody elicitation. *Cell Rep*. 2017;19:719–32. doi:10.1016/j.celrep.2017.04.013.
39. Wang H, Gristick HB, Scharf L, West AP, Galimidi RP, Seaman MS, Freund NT, Nussenzweig MC, Bjorkman PJ. Asymmetric recognition of HIV-1 Envelope trimer by V1V2 loop-targeting antibodies. *eLife*. 2017;6:e27389. doi:10.7554/eLife.27389.
40. Medina-Ramírez M, Garces F, Escolano A, Skog P, SW T, ID M-S, AT M, Yasmeen A, A-J B, Ozorowski G, et al. Design and crystal structure of a native-like HIV-1 envelope trimer that engages multiple broadly neutralizing antibody precursors in vivo. *J Exp Med*. 2017;214:2573–90. doi:10.1084/jem.20161160.
41. Sarkar A, Bale S, Behrens A-J, Kumar S, Sharma SK, Val N, Pallesen J, Irimia A, Diwanji DC, Stanfield RL, et al. Structure of a cleavage-independent HIV Env recapitulates the glycoprotein architecture of the native cleaved trimer. *Nat Commun*. 2018;9:1–14. doi:10.1038/s41467-018-04272-y.
42. Moyo T, Ereño-Orbea J, Jacob RA, Pavillet CE, Kariuki SM, Tangie EN, Julien J-P, Dorfman JR. Molecular basis of unusually high neutralization resistance in tier 3 HIV-1 strain 253-11. *J Virol*. 2018;92:e02261–17. doi:10.1128/JVI.02261-17.
43. Xu K, Acharya P, Kong R, Cheng C, Chuang G-Y, Liu K, Louder MK, O'Dell S, Rawi R, Sastry M, et al. Epitope-based vaccine design yields fusion peptide-directed antibodies that neutralize diverse strains of HIV-1. *Nat Med*. 2018;24:857–67. doi:10.1038/s41591-018-0042-6.
44. Barnes CO, Gristick HB, Freund NT, Escolano A, Lyubimov AY, Hartweger H, West AP, Cohen AE, Nussenzweig MC, Bjorkman PJ. Structural characterization of a highly-potent V3-glycan broadly neutralizing antibody bound to natively-glycosylated HIV-1 envelope. *Nat Commun*. 2018;9:1–12. doi:10.1038/s41467-018-03632-y.
45. Zhang P, Gorman J, Geng H, Liu Q, Lin Y, Tsybovsky Y, Go EP, Dey B, Andine T, Kwon A, et al. Interdomain stabilization impairs CD4 binding and improves immunogenicity of the HIV-1 envelope trimer. *Cell Host Microbe*. 2018;23:832–844.e6. doi:10.1016/j.chom.2018.05.002.
46. Bhiman JN, Anthony C, Doria-Rose NA, Karimanzira O, Schramm CA, Khoza T, Kitchin D, Botha G, Gorman J, Garrett NJ, et al. Viral variants that initiate and drive maturation of V1V2-directed HIV-1 broadly neutralizing antibodies. *Nat Med*. 2015;21:1332–36. doi:10.1038/nm.3963.
47. Zhou T, Lynch RM, Chen L, Acharya P, Wu X, Doria-Rose NA, Joyce MG, Lingwood D, Soto C, Bailer RT, et al. Structural Repertoire of HIV-1-Neutralizing Antibodies Targeting the CD4 Supersite in 14 Donors. *Cell*. 2015;161:1280–92. doi:10.1016/j.cell.2015.05.007.
48. Liao H-X, Lynch R, Zhou T, Gao F, Alam SM, Boyd SD, Fire AZ, Roskin KM, Schramm CA, Zhang Z, et al. Co-evolution of a broadly neutralizing HIV-1 antibody and founder virus. *Nature*. 2013;496:469–76. doi:10.1038/nature12053.
49. Schanz M, Liechti T, Zagordi O, Miho E, Reddy ST, Günthard HF, Trkola A, Huber M. High-throughput sequencing of human immunoglobulin variable regions with subtype identification. *Plos One*. 2014;9:e111726. doi:10.1371/journal.pone.0111726.
50. Soto C, Ofek G, Joyce MG, Zhang B, McKee K, Longo NS, Yang Y, Huang J, Parks R, Eudailey J, et al. Developmental pathway of the MPER-directed HIV-1-neutralizing antibody 10E8. *Plos One*. 2016;11:e0157409. doi:10.1371/journal.pone.0157409.
51. Zhu J, Wu X, Zhang B, McKee K, O'Dell S, Soto C, Zhou T, Casazza JP, Mullikin JC, Kwong PD, et al. De novo identification of VRC01 class HIV-1-neutralizing antibodies by next-generation sequencing of B-cell transcripts. *Proc Natl Acad Sci U S A*. 2013;110:E4088–E4097. doi:10.1073/pnas.1306262110.
52. Kepler TB, Liao H-X, Alam SM, Bhaskarabhatla R, Zhang R, Yandava C, Stewart S, Anasti K, Kelsoe G, Parks R, et al. Immunoglobulin gene insertions and deletions in the affinity maturation of HIV-1 broadly reactive neutralizing antibodies. *Cell Host Microbe*. 2014;16:304–13. doi:10.1016/j.chom.2014.08.006.
53. Ester M, Kriegel H-P, Sander J, Xu X A density-based algorithm for discovering clusters a density-based algorithm for discovering clusters in large spatial databases with noise. in *Proceedings of the Second International Conference on Knowledge Discovery and Data Mining KDD'96*. (Portland, Oregon: AAAI Press), 226–31. [Accessed February 1, 2017]. <http://dl.acm.org/citation.cfm?id=3001460.3001507>
54. Wall JS, Gupta V, Wilkerson M, Schell M, Loris R, Adams P, Solomon A, Stevens F, Dealwis C. Structural basis of light chain amyloidogenicity: comparison of the thermodynamic properties, fibrillogenic potential and tertiary structural features of four VL6 proteins. *J Mol Recognit*. 2004;17:323–31. doi:10.1002/jmr.681.
55. Hayward S. Peptide-plane flipping in proteins. *Protein Sci*. 2001;10:2219–27. doi:10.1110/ps.23101.
56. North B, Lehmann A, Dunbrack RL. A new clustering of antibody CDR loop conformations. *J Mol Biol*. 2011;406:228–56.
57. Galson JD, Clutterbuck EA, Trück J, Ramasamy MN, Münz M, Fowler A, Cerundolo V, Pollard AJ, Lunter G, Kelly DF. BCR repertoire sequencing: different patterns of B cell activation after two Meningococcal vaccines. *Immunol Cell Biol*. 2015;93:885–95. doi:10.1038/icb.2015.57.
58. Ellebedy AH, Jackson KJL, Kissick HT, Nakaya HI, Davis CW, Roskin KM, McElroy AK, Oshansky CM, Elbein R, Thomas S, et al. Defining antigen-specific plasmablast and memory B cell subsets in human blood after viral infection or vaccination. *Nat Immunol*. 2016;17:1226–34. doi:10.1038/ni.3533.
59. Johnson EL, Doria-Rose NA, Gorman J, Bhiman JN, Schramm CA, Vu AQ, Law WH, Zhang B, Bekker V, Abdool Karim SS, et al. Sequencing HIV-neutralizing antibody exons and introns reveals detailed aspects of lineage maturation. *Nat Commun*. 2018;9:1–13. doi:10.1038/s41467-018-06424-6.
60. Rubelt F, Bolen CR, McGuire HM, Heiden JAV, Gadala-Maria D, Levin M, Euskirchen G, Mamedov MR, Swan GE, Dekker CL, et al. Individual heritable differences result in unique cell lymphocyte receptor repertoires of naïve and antigen-experienced cells. *Nat Commun*. 2016;7:1–12. doi:10.1038/ncomms11112.
61. Turchaninova MA, Davydov A, Britanova OV, Shugay M, Bikos V, Egorov ES, Kirgizova VI, Merzlyak EM, Staroverov DB, Bolotin DA, et al. High-quality full-length immunoglobulin profiling with unique molecular barcoding. *Nat Protoc*. 2016;11:1599–616. doi:10.1038/nprot.2016.093.

62. Vander Heiden JA, Stathopoulos P, Zhou JQ, Chen L, Gilbert TJ, Bolen CR, Barohn RJ, Dimachkie MM, Ciafaloni E, Broering TJ, et al. Dysregulation of B cell repertoire formation in myasthenia gravis patients revealed through deep sequencing. *J Immunol Baltim Md* 1950. 2017;198:1460–73. doi:10.4049/jimmunol.1601415.
63. Kirik U, Persson H, Levander F, Greiff L, Ohlin M. Antibody heavy chain variable domains of different germline gene origins diversify through different paths. *Front Immunol*. 2017;8. doi:10.3389/fimmu.2017.01433.
64. Modi V, Dunbrack RL. Defining a new nomenclature for the structures of active and inactive kinases. *Proc Natl Acad Sci*. 2019;116:6818–27. doi:10.1073/pnas.1814279116.
65. Shapovalov M, Vucetic S, Dunbrack RL. A new clustering and nomenclature for beta turns derived from high-resolution protein structures. *PLoS Comput Biol*. 2019;15:e1006844. doi:10.1371/journal.pcbi.1006844.
66. Adolf-Bryfogle J, Kalyuzhnyi O, Kubitz M, Weitzner BD, Hu X, Adachi Y, Schief WR, Dunbrack RL. Rosetta antibody design (RABD): a general framework for computational antibody design. *bioRxiv*. 2018:183350. doi:10.1101/183350.
67. Weitzner BD, Jeliakov JR, Lyskov S, Marze N, Kuroda D, Frick R, Adolf-Bryfogle J, Biswas N, Dunbrack RL, Gray JJ. Modeling and docking antibody structures with Rosetta. *Nat Protoc*. 2017;12:401–16. doi:10.1038/nprot.2016.180.
68. Dunbar J, Krawczyk K, Leem J, Marks C, Nowak J, Regep C, Georges G, Kelm S, Popovic B, Deane CM. SABPred: a structure-based antibody prediction server. *Nucleic Acids Res*. 2016;44:W474–W478. doi:10.1093/nar/gkw361.
69. Marcatili P, Rosi A, Tramontano A. PIGS: automatic prediction of antibody structures. *Bioinformatics*. 2008;24:1953–54. doi:10.1093/bioinformatics/btn341.
70. Adolf-Bryfogle J, Xu Q, North B, Lehmann A, Dunbrack RL. PyIgClassify: a database of antibody CDR structural classifications. *Nucleic Acids Res*. 2015;43:D432–D438.
71. Honegger A, Plückthun AP. Yet Another Numbering Scheme for Immunoglobulin Variable Domains: An Automatic Modeling and Analysis Tool. *Journal of Molecular Biology*. 2001;309:657–670 3 doi:10.1006/jmbi.2001.4662
72. Leaver-Fay A, Tyka M, Lewis SM, Lange OF, Thompson J, Jacak R, Kaufman K, Renfrew PD, Smith CA, Sheffler W, et al. ROSETTA3: an object-oriented software suite for the simulation and design of macromolecules. *Methods Enzymol*. 2011;487:545–74. doi:10.1016/B978-0-12-381270-4.00019-6.
73. Shrake A, Rupley JA. Environment and exposure to solvent of protein atoms. Lysozyme and insulin. *J Mol Biol*. 1973;79:351–71. doi:10.1016/0022-2836(73)90011-9.
74. Pedregosa F, Varoquaux G, Gramfort A, Michel V, Thirion B, Grisel O, Blondel M, Müller A, Nothman J, Louppe G, et al. Scikit-learn: machine Learning in Python. *ArXiv E-Prints*. 2012;1201:arXiv:1201.0490.
75. Simpson GG. Notes on the measurement of faunal resemblance. *Am J Sci*. 1960;258:300–11.
76. Kovaltsuk A, Leem J, Kelm S, Snowden J, Deane CM, Krawczyk K. Observed antibody space: a resource for data mining next-generation sequencing of antibody repertoires. *J Immunol*. 2018;201:2502–09. doi:10.4049/jimmunol.1800708.
77. Stewart JJ, Lee CY, Ibrahim S, Watts P, Shlomchik M, Weigert M, Litwin SA. Shannon entropy analysis of immunoglobulin and T cell receptor. *Mol Immunol*. 1997;34:1067–82. doi:10.1016/S0161-5890(97)00130-2.
78. Giudicelli V, Brochet X, Lefranc M-P. IMGT/V-QUEST: IMGT Standardized Analysis of the Immunoglobulin (IG) and T Cell Receptor (TR) Nucleotide Sequences. *Cold Spring Harb Protoc*. 2011;2011:pdb.prot5633. doi:10.1101/pdb.prot5633.