

## **RAIN: Machine Learning-based identification for HIV-1 bNAbs**

Mathilde Foglierini<sup>1,2,3,\*</sup>, Pauline Nortier<sup>1,2,\*</sup>, Rachel Schelling<sup>1,2</sup>, Rahel R. Winiger<sup>1,2</sup>, Philippe Jacquet<sup>4</sup>, Sijy O'Dell<sup>5</sup>, Davide Demurtas<sup>6</sup>, Maxmillian Mpina<sup>7</sup>, Omar Lweno<sup>7</sup>, Yannick D. Muller<sup>1,2</sup>, Constantinos Petrovass<sup>8</sup>, Claudia Daubenberger<sup>9,10</sup>, Mathieu Perreau<sup>1</sup>, Nicole A Doria-Rose<sup>5</sup>, Raphael Gottardo<sup>3</sup> and Laurent Perez<sup>1,2,#</sup>

<sup>1</sup>Department of Medicine, Service of Immunology and Allergy, Lausanne University Hospital and University of Lausanne, Lausanne, Switzerland.

<sup>2</sup>Centre for Human Immunology Lausanne, Switzerland.

<sup>3</sup>Biomedical Data Science Centre, Lausanne University Hospital and University of Lausanne, Lausanne, Switzerland.

<sup>4</sup>Scientific Computing and Research Support Unit, University of Lausanne, Lausanne, Switzerland.

<sup>5</sup>Vaccine Research Center, National Institute of Allergy and Infectious Diseases, National Institutes of Health, Bethesda, MD, USA.

<sup>6</sup>Interdisciplinary center of electron microscopy, CIME, Ecole Polytechnique Fédérale de Lausanne, Lausanne, Switzerland.

<sup>7</sup>Ifakara Health Institute, Bagamoyo, United Republic of Tanzania.

<sup>8</sup>Department of Laboratory Medicine and Pathology, Institute of Pathology, Lausanne University Hospital, Lausanne, Switzerland

<sup>9</sup>Department of Medical Parasitology and Infection Biology, Clinical Immunology Unit, Swiss Tropical and Public Health Institute, Basel, Switzerland.

<sup>10</sup>University of Basel, Basel, Switzerland.

\*These authors contributed equally

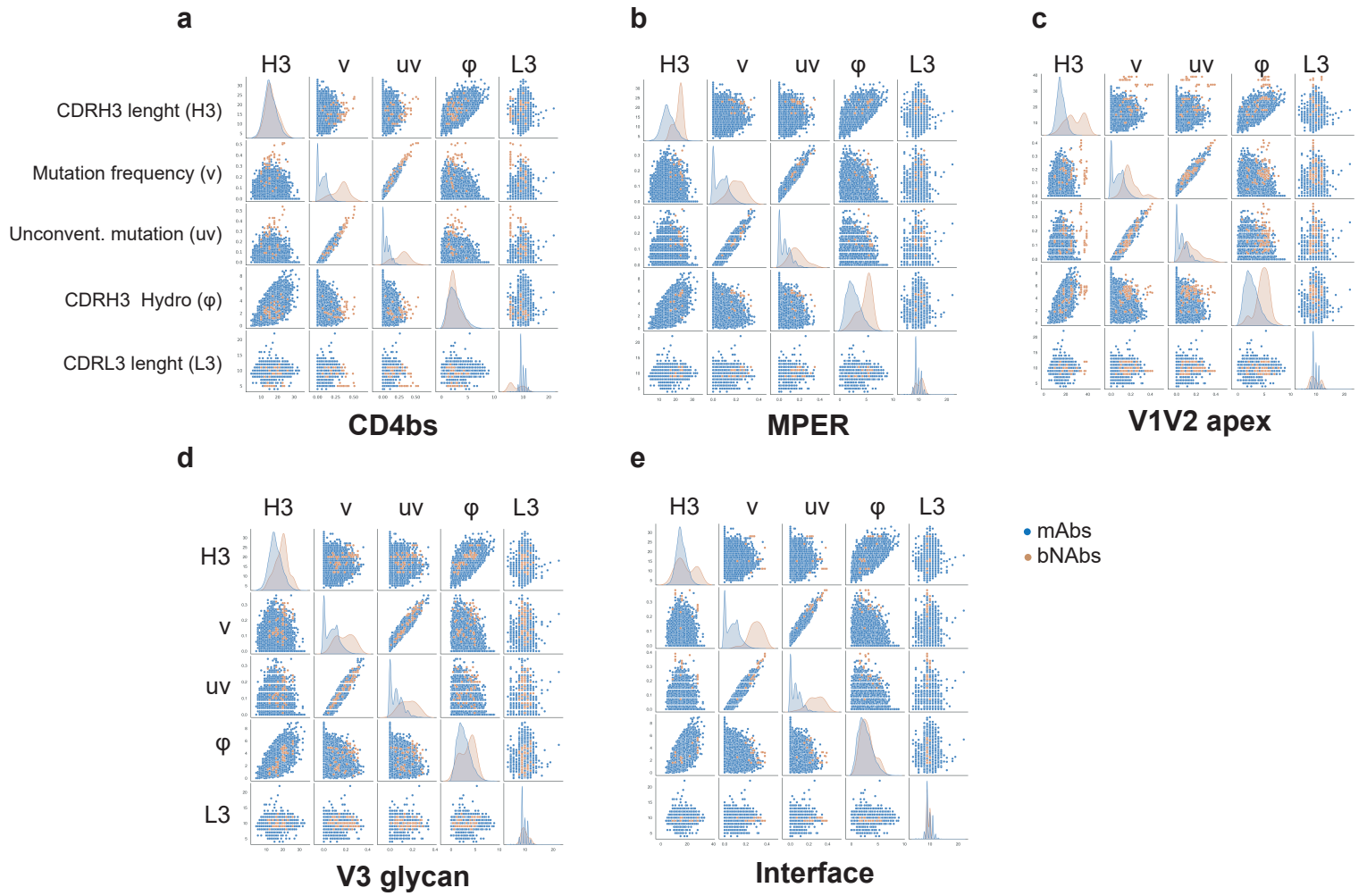
#Corresponding author: Laurent Perez

**email:** [laurent.perez@chuv.ch](mailto:laurent.perez@chuv.ch)

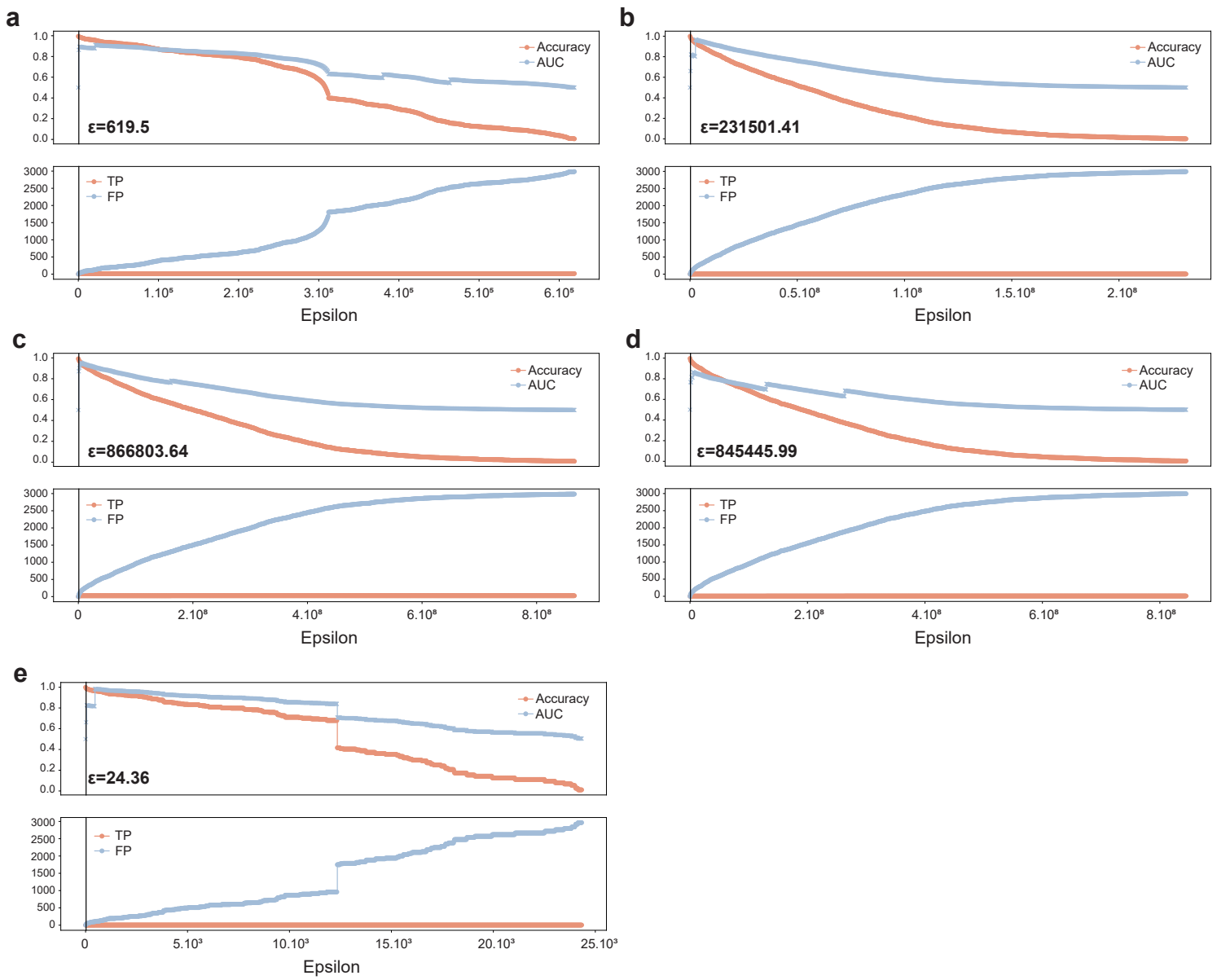
## **Supplementary Information**

Supplementary Figures 1-10

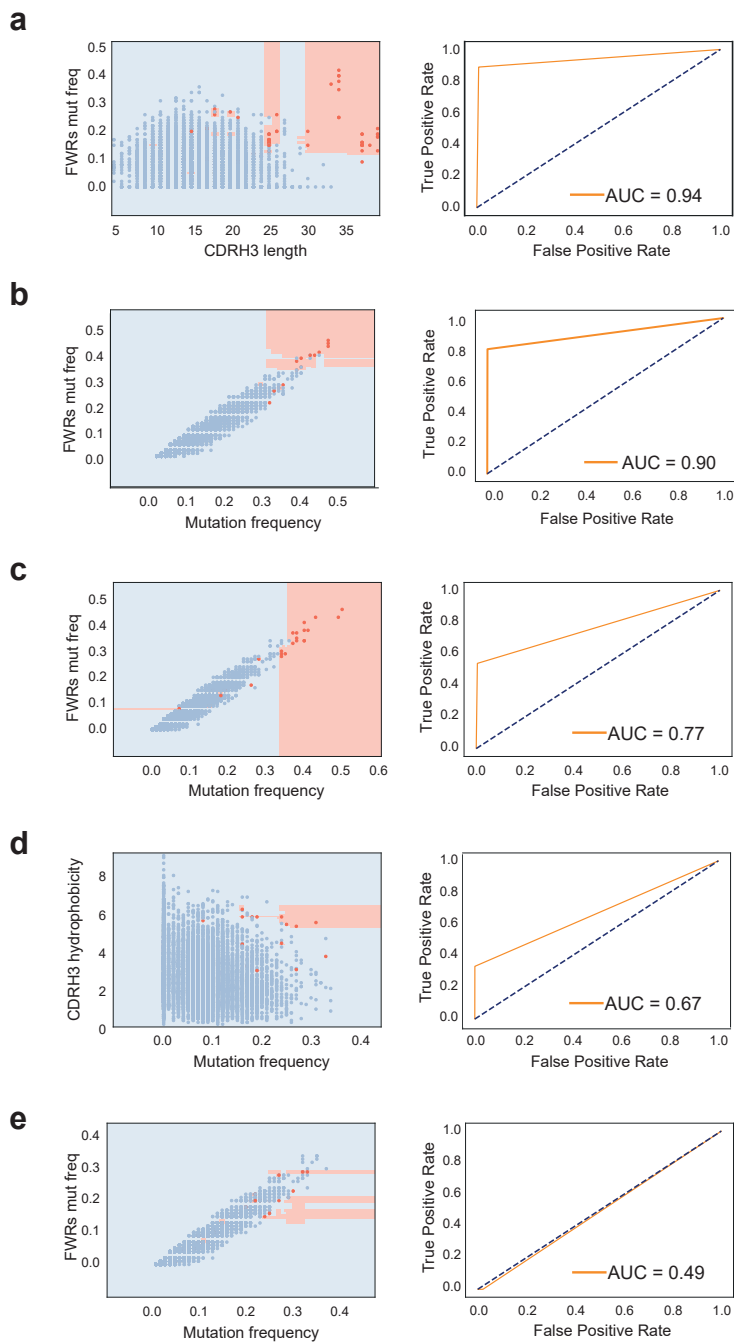
Supplementary Tables 1-3



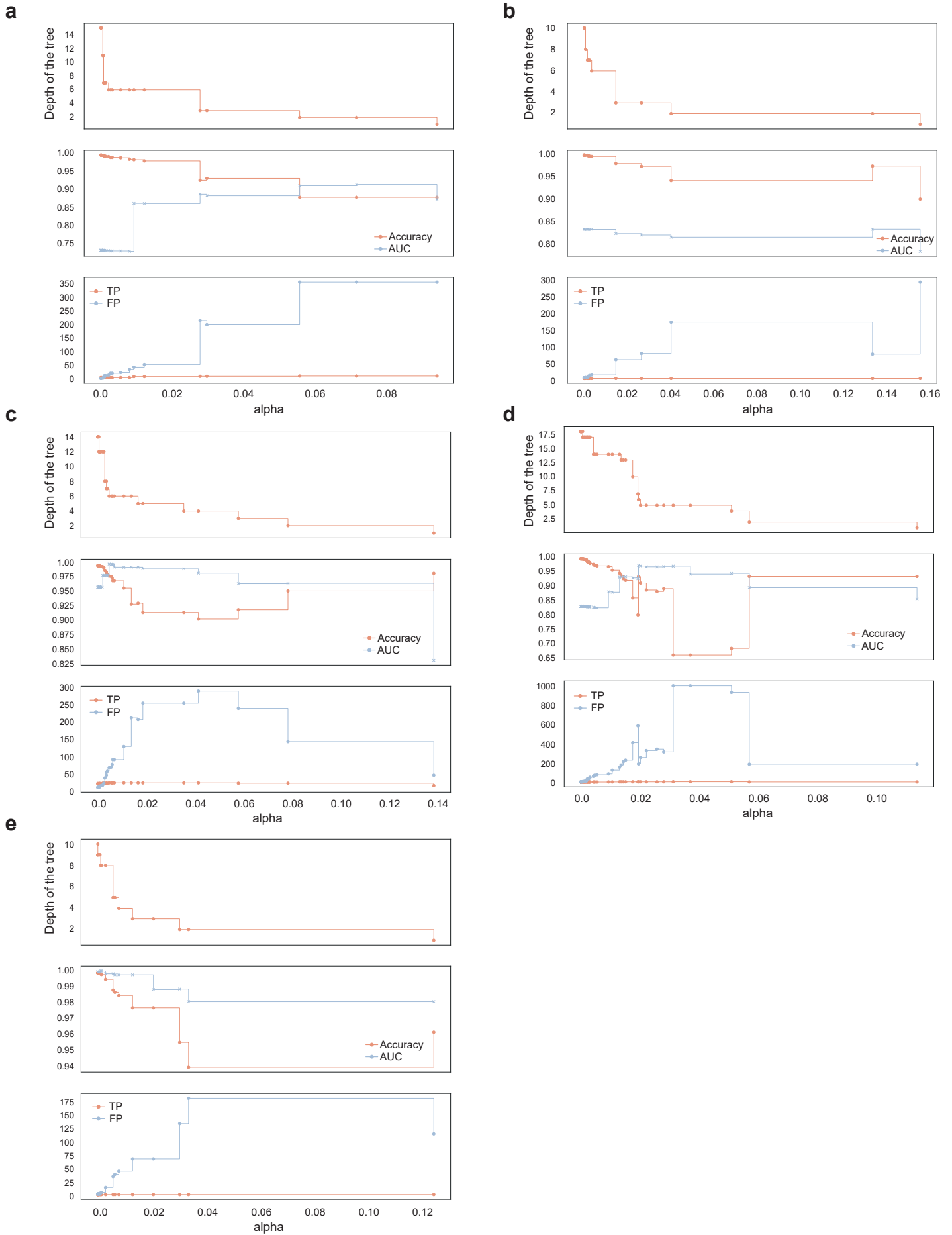
**Supplementary Figure. 1. Correlation of each feature pair in the dataset against the class distribution.** Pairwise relationships between the specific features comparing mAbs (blue) and bNAbs (red). Each bNAb category is represented by a single plot per antigenic site: **(a)** CD4bs, **(b)** MPER, **(c)** V1V2 apex, **(d)** V3 glycan, and **(e)** gp120/gp41 interface. Source data are provided as a Source Data file.



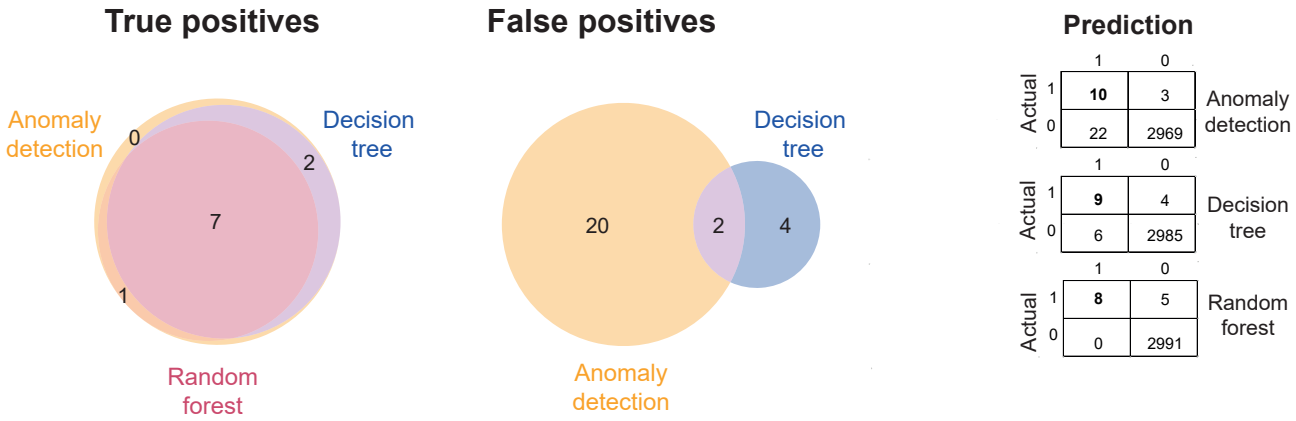
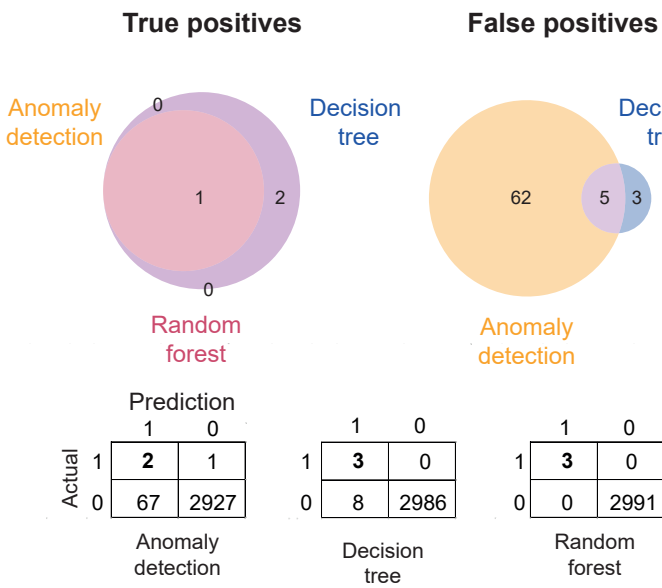
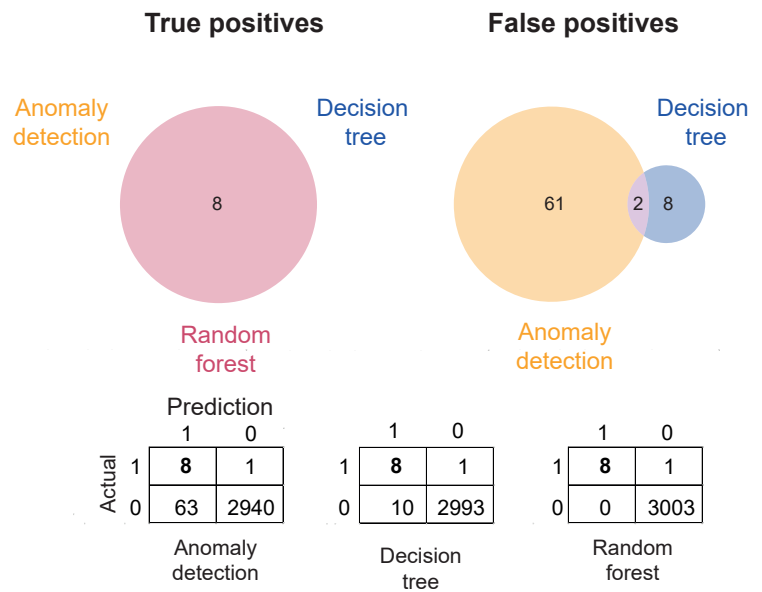
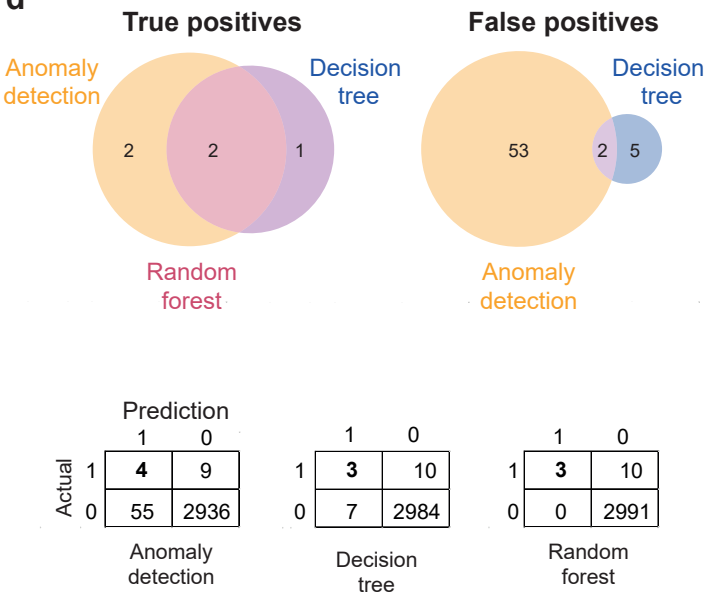
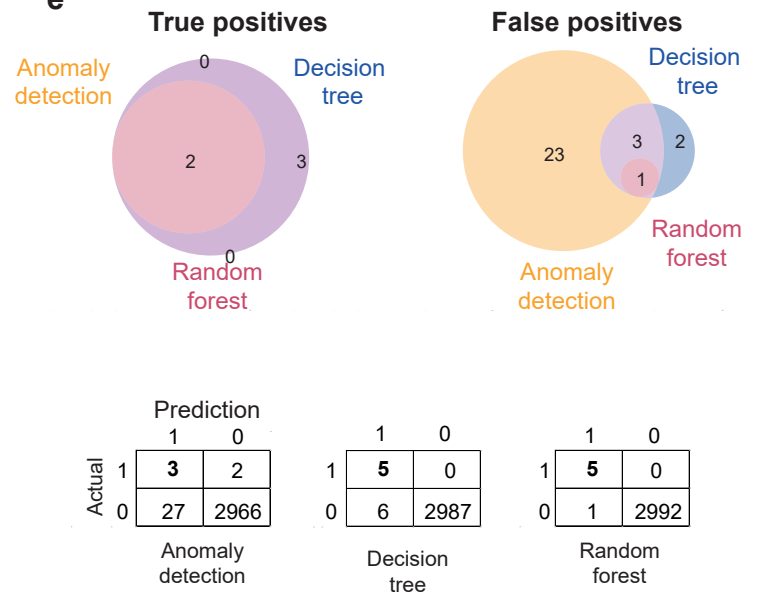
**Supplementary Figure 2. Determination of the epsilon parameter for the Anomaly detection algorithm.** Epsilon parameters for the antigenic site (a) CD4bs, (b) MPER, (c) V1V2 apex, (d) V3 glycan, and (e) gp120/gp41 interface. The accuracy is in red and the area under the curve (AUC) is in blue.



**Supplementary Figure 3. Random Forest classifier with only two features.** Left panels are the decision regions obtained from the training dataset and right panels are the receiver-operating characteristic (ROC) curves in orange with the corresponding area under the curve (AUC). The dashed blue line is for visual reference. Statistics are visualized for the different antigenic sites: **(a)** V1V2 apex, **(b)** gp120/gp41 interface, **(c)** CD4bs, **(d)** MPER, and **(e)** V3 glycan.



**Supplementary Figure 4. Determination of the cost complexity pruning parameter  $\alpha$  ( $ccp\_alpha$ ) for the Decision Tree classifier. Targeted antigenic sites are (a) CD4bs, (b) MPER, (c) V1V2 apex, (d) V3 glycan, and (e) gp120/gp41 interface.**

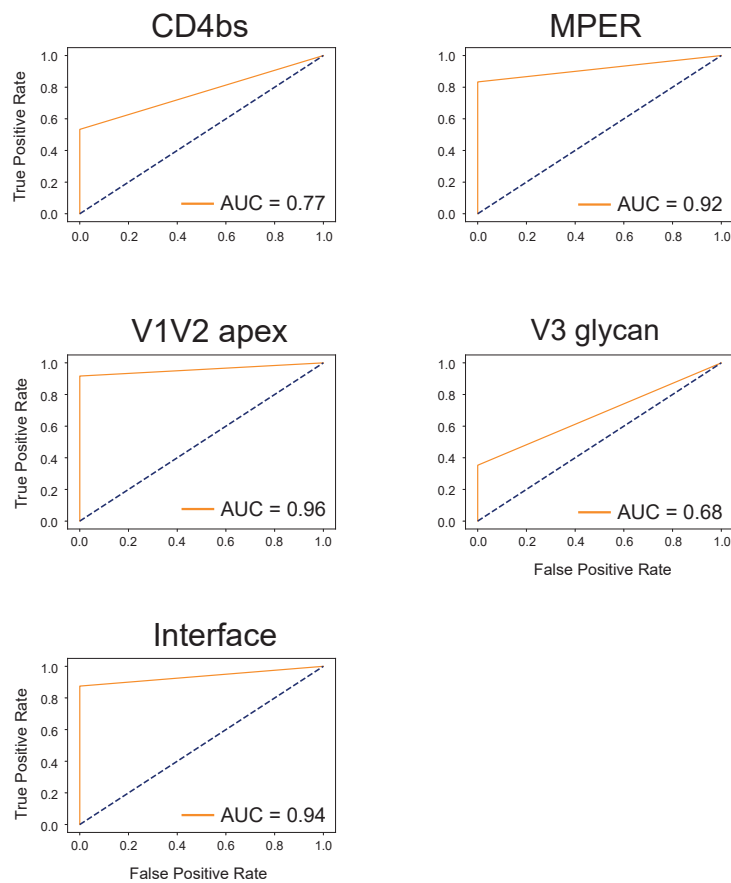
**a****b****c****d****e**

**Supplementary Figure 5. Overlap of false and true positives between the three different models using the test datasets.** Venn diagrams and Confusion matrices for each model are represented next to the Venn diagram. Antigenic sites are **(a)** CD4bs, **(b)** MPER, **(c)** V1V2 apex, **(d)** V3 glycan, and **(e)** gp120/gp41 interface. The yellow circle represents AD, the blue one is DT, and the red one is RF. Source data are provided as a Source Data file.

**a**

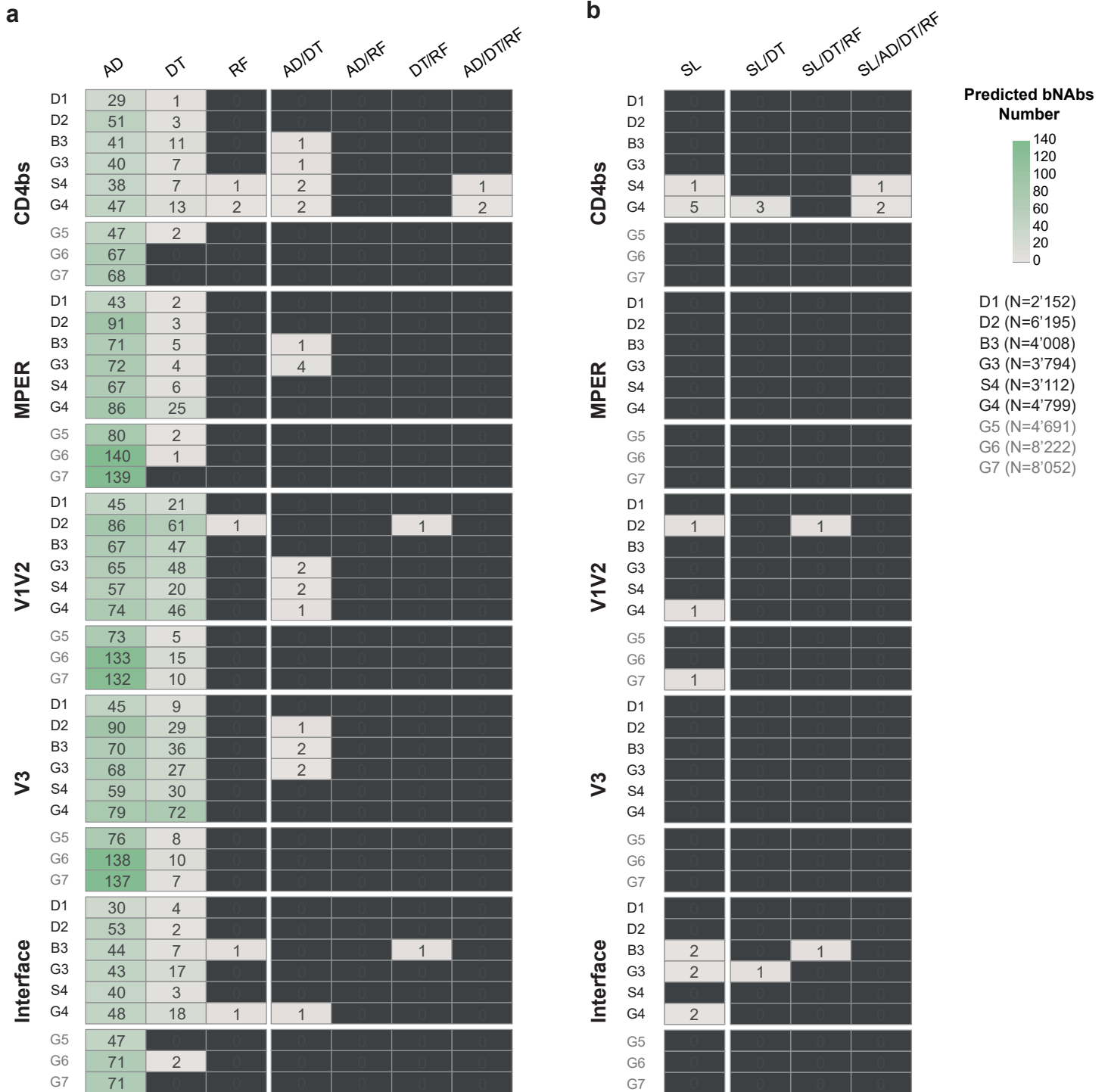
	TP	FP	TN	FN	AUC	Accuracy	Recall	Precision
<b>CD4bs</b>	8	0	3739	7	0.77	1.00	0.53	1.00
<b>MPER</b>	5	0	3740	1	0.92	1.00	0.83	1.00
<b>V1V2</b>	11	0	3753	1	0.96	1.00	0.92	1.00
<b>V3</b>	6	0	3738	11	0.68	1.00	0.35	1.00
<b>Interface</b>	7	0	3739	1	0.94	1.00	0.88	1.00

score  
1  
0.8  
0.6  
0.4  
0.2  
0

**b****Supplementary Figure 6. Performance and results of the Super Learner Ensembles algorithm.**

**(a)** Performance metrics of the algorithm using the test dataset with Accuracy =  $(TP+TN) / (TP+FP+TN+FN)$ , Recall =  $TP / (TP+FN)$ , and Precision =  $TP / (TP+FP)$ . **(b)** Receiver-operating characteristic (ROC) curves and corresponding area under the curve (AUC) statistics for each bNAb antigenic site with the test dataset.





**Supplementary Figure 7. Predicted bNAb per dataset.** (a) The heatmap illustrates the number of predicted bNAb for each antigenic site and each run of donors using various algorithms, with D1: donor 1, D2: donor 2 (donor that did not have sera with broad neutralization activity) and B3, G3, S4, G4: donor 3 (donor serum with broad neutralization activity). The G5, G6, and G7 datasets correspond to a Influenza-specific repertoire. The columns on the right side of the heatmap show bNAb identified by different algorithms. The abbreviations used in each cell represent the combination of algorithms that share the specific bNAb (AD/DT: bNAb predicted only by both AD and DT, AD/RF: bNAb predicted only by both AD and RF, DT/RF: bNAb predicted only by both DT and RF and AD/DT/RF: bNAb predicted by all AD, DT and RF). (b) The heatmap illustrates the number of predicted bNAb for each antigenic site and each run of donors using the Super Learner (SL) algorithm. The columns on the right side of the heatmap show bNAb shared by different algorithms. The abbreviations used in each cell represent the combination of algorithms that share the specific bNAb (SL/DT: bNAb predicted only by both SL and DT, SL/DT/RF: bNAb predicted only by SL, DT and RF and SL/AD/DT/RF: bNAb predicted by all SL, AD, DT and RF).

**a**

		FR1	CDR1	FR2	CDR2	FR3	CDR3	FR4																										
UCA	1	5	10	15	20	25	30	35	40	45	50	52A	55	60	65	70	75	80	82A	82B	82C	85	90	95	100	100A	100B	100C	100D	100E	100F	105	110	113
UCA	QVQLVQSGAEVKKPGASVKVSKASGYTFTGYMHWRQAPGQGLEWMGWINPNSGGTNYAQKFQGRVTMTRDTSISTAYMELSLRSLRSDDTAVYYCAXXXXXXXXXXXXXXXXXXWQGGLVTVSS																																	
VRC01	.....GQM.....E.MRI..R.....E.IDCTLN..I.L...KRP.....LK.RG.AV...RPL.....VYSD..FL..RS.TV.....F.TRGKNCDYNWDF--EH...R..P.I...																																	
bNAb4251	..HVM...DQ.....Q...TT..SS.IEDSL..IQ.V...EP..L..V..RH.AV..SW.IRD.I.....I.K..M.VQMRG.Q.....M...KSRRGANWA---L...W..RI....																																	
bNAb2101	..VM...DQ.RE.....R...T.ED.VES.L.....EP..LA.....RN.AV...-SLRD.L.L...IY...V.VDMRG.Q.....KARRGNTWAF---R...W..RI....																																	
bNAb1586	.E.P...PSL.....T...RGDENL.IE..I..I.....H...MSILT.AP..SGN.RN.MS.Y..R.....D.RG.T.....V..TSRRRSGRGGTWFQ.....																																	

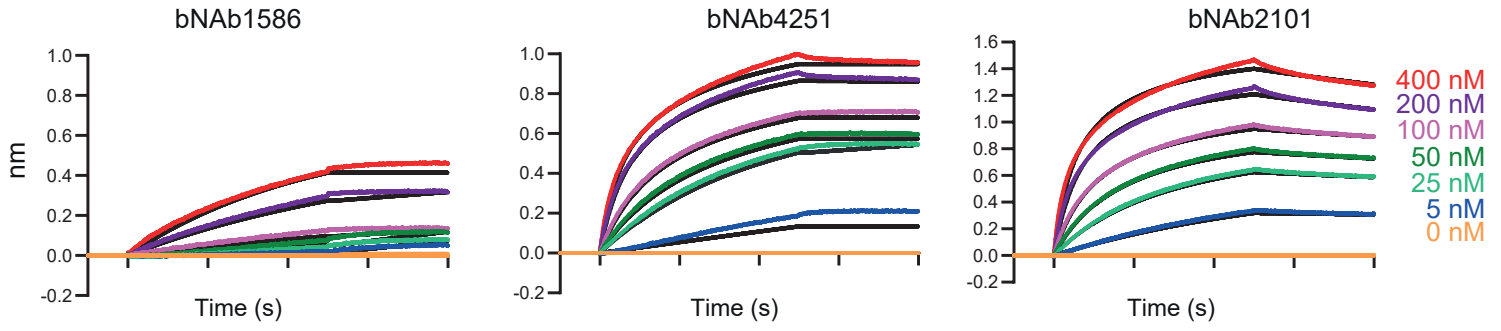
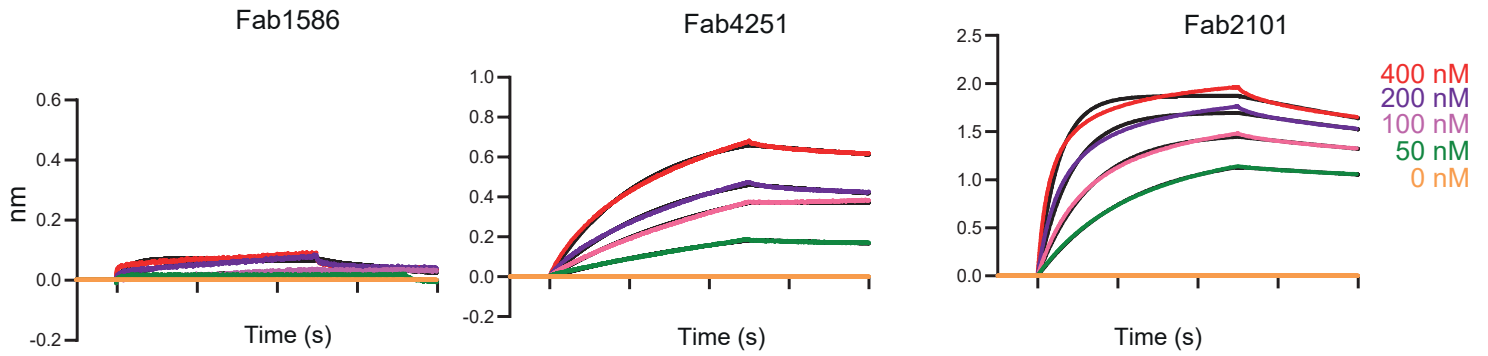
**b**

		FR1	CDR1	FR2	CDR2	FR3	CDR3	FR4															
UCA	1	5	10	15	20	25	30	35	40	45	50	55	60	65	70	75	80	85	90	95	100	105	107
UCA	DIQMTQSPSSLSASVGDRTITTCQASQDISNYLNWYQQKPKGAPKLLIYDASNLETGVPFRFSGSGSGTDFTFTTISLQPEDIAITYYCAQXX---XXFGGGTKVEIK																						
bNAb4251	...L...Y.A.....RE.N---D...L.R...P.....SG..K..R.....R..S...-SL...G.....G.....VF...EF.....R.D..																						
bNAb2101	.L.LD.....R..D---R.....P.....K.DR.....A...-L..NT.E.D.F.....VF...QF.....D..																						
bNAb1586	.T.V...P.....GT.....P..G.Q.Q.H.....R.....G.R..R.....R.SL..NN.....F.....AF...YS.....R....																						

**c**

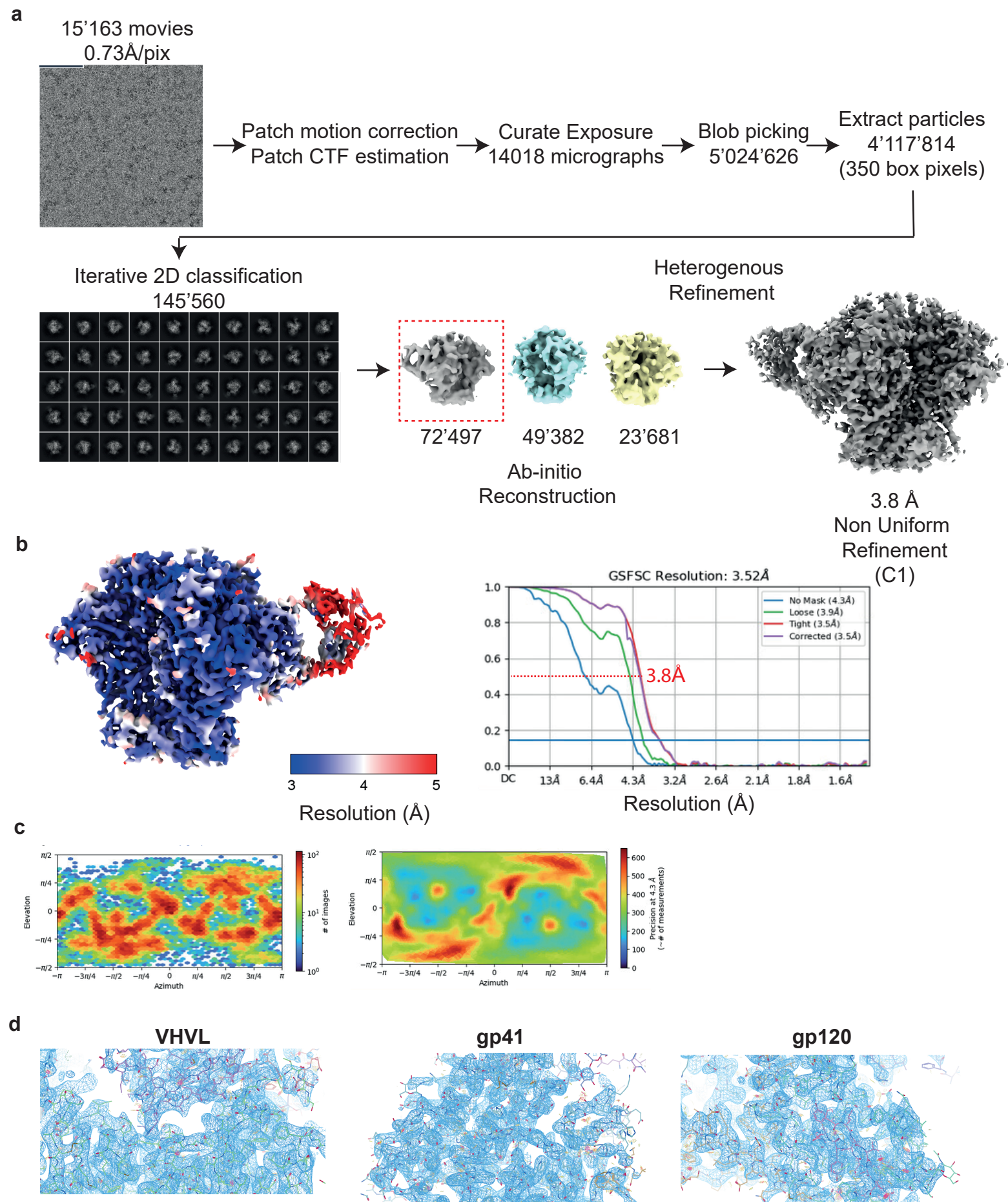
		FR1	CDR1	FR2	CDR2	FR3	CDR3	FR4															
UCA	1	5	10	15	20	25	30	35	40	45	50	55	60	65	70	75	80	85	90	95	100	105	107
UCA	EIVLTQSPATLSLSPGERATLSCRASQSVSSYLAWYQQKPGQAPRLLIYDASNRTATGIPARFSGSGSGTDFTLTISSLEPEDFAVYYCAQXX---XXFGGGTKVQVD																						
VRC01	.....G.....T..I.I...T...-YGS.....R.....V..SG.T..A...D.....RW.P.YN...N..SG..G.....Y...EF.....																						

**Supplementary Figure 8. bNAb alignment of VH and VK amino acid sequences.** Upper panel: (a) Alignment of the VH sequence of the three predicted bNAb with their UCA (unmutated common ancestor) and VRC01 bNAb. UCA is constituted of VH1-2\*02 and JH4\*02, the D gene is masked and represented by X. Residue positions are according to Kabat numbering. Dots indicate identical residues. (b) Middle panel: alignment of the VK sequences of the three predicted bNAb with their UCA (unmutated common ancestor). UCA is constituted of VK1-33\*01 and JK4\*01, the LCDR3 is masked and represented by X. (c) Lower panel: alignment of the VK sequence of the VRC01 with its UCA VK3-11\*01 and JK2\*01.

**a****b****Supplementary Figure 9. Antibodies and Fabs interaction with BG505-DS-SOSIP.**

**(a)** Bivalent analyte fitting for antibody SOSIP interactions with concentrations ranging from 5 to 400 nM.

**(b)** Fitting of Fab interaction with BG505-DS-SOSIP using 1:1 model. For Fabs, concentrations ranging from 50 to 400 nM were used. Source data are provided as a Source Data file.



**Supplementary Figure 10. Cryo-EM data processing and validation of Fab-BG505-DS-SOSIP complexes.**

(a) Representative cryo-EM micrograph, 2D class average images and data processing flow chart.

(b) Local resolution of the final map and Gold standard Fourier shell correlation (FSC) at 0.143 resolution of 3.8 Å.

(c) Angular distribution of the Fab-BG505-DS-SOSIP particles in the final round of 3D refinement.

(d) Different density maps shown at threshold of  $6\sigma$  for VHVL, gp41 and gp120, respectively.

**Supplementary Table 1. Source of the ten datasets used as mAbs in the machine learning models and the three datasets used as Influenza repertoires.**

	<b>Paired BCRs (IGK/L+IGH)</b>	<b>Reference</b>
<b>H1</b>	4 581	ArrayExpress: E-MTAB-11174 (Memory B cells of donor 1)
<b>H2</b>	4 740	ArrayExpress: E-MTAB-11174 (Memory B cells of donor 2)
<b>H3</b>	965	SRA: SRR17717616
<b>H4</b>	903	SRA: SRR17717597
<b>H5</b>	895	SRA: SRR17717612
<b>H6</b>	413	SRA: SRR17717605
<b>H7</b>	341	SRA: SRR17717601
<b>H8</b>	786	SRA: SRR17717593
<b>H9</b>	429	10X dataset: 10k Human PBMCs, 5' v2.0, Chromium Controller, Single Cell Immune Profiling Dataset by Cell Ranger 6.1.0 (2021, August 9)
<b>H10</b>	909	10X dataset: Human PBMC from a Healthy Donor, 10k cells (v2), Single Cell Immune Profiling Dataset by Cell Ranger 5.0.0 (2020, November 19)
<b>G5</b>	4 691	BioSample SAMN07733010, Day 9 after influenza vaccination. SRA : SRR10596386, SRR10596375, SRR10596364, SRR10596353
<b>G6</b>	8 222	BioSample SAMN07733010, Day 7 after influenza vaccination. SRA : SRR10596411, SRR10596412, SRR10596410, SRR10596409
<b>G7</b>	8 052	BioSample SAMN07733010, Day 7 after influenza vaccination. SRA : SRR10596407, SRR10596406, SRR10596405, SRR10596404

**Supplementary Table 2. Performance metrics of the three algorithms using the validation dataset.**

<b>Algo</b>	<b>Ag site</b>	<b>bNAbs Numbers</b>	<b>TP</b>	<b>FP</b>	<b>TN</b>	<b>FN</b>	<b>AUC</b>	<b>Accuracy</b>	<b>Recall</b>	<b>Precision</b>
<b>AD</b>	CD4bs	15	12	22	2966	3	0.90	0.99	0.80	0.35
<b>DT</b>	CD4bs	15	7	4	2984	8	0.73	1.00	0.47	0.64
<b>RF</b>	CD4bs	15	7	0	2988	8	0.90	1.00	0.47	1.00
<b>AD</b>	MPER	3	1	38	2956	2	0.66	0.99	0.33	0.03
<b>DT</b>	MPER	3	2	4	2990	1	0.83	1.00	0.67	0.33
<b>RF</b>	MPER	3	2	0	2994	1	0.83	1.00	0.67	1.00
<b>AD</b>	V1V2 apex	25	19	55	2932	6	0.87	0.98	0.76	0.26
<b>DT</b>	V1V2 apex	25	23	12	2975	2	0.96	1.00	0.92	0.66
<b>RF</b>	V1V2 apex	25	19	0	2987	6	1.00	1.00	0.76	1.00
<b>AD</b>	V3 glycan	9	5	57	2938	4	0.77	0.98	0.56	0.08
<b>DT</b>	V3 glycan	9	6	8	2987	3	0.83	1.00	0.67	0.43
<b>RF</b>	V3 glycan	9	2	0	2995	7	1.00	1.00	0.22	1.00
<b>AD</b>	Interface	3	1	31	2964	2	0.66	0.99	0.33	0.03
<b>DT</b>	Interface	3	3	4	2991	0	1.00	1.00	1.00	0.43
<b>RF</b>	Interface	3	2	0	2995	1	1.00	1.00	0.67	1.00

### Supplementary Table 3. Cryo-EM data collection, refinement, and validation statistics.

Data Accession PDB ID: 8S2E EMDB: EMD-19665

---

Magnification	165'000
Voltage (kV)	300
Electron exposure (e-/Å <sup>2</sup> )	39.89
Defocus range (μm)	-0.9 to -2.4
Pixel size (Å)	0.73277
Symmetry imposed	C1
Micrographs collected (no.)	15,163
Final particle images (no.)	72497
Map resolution (Å)	3.8
FSC threshold	0.143
Map resolution range (Å)	
<b>Refinement</b>	
Initial model used (PDB code)	4TVP
Model resolution (Å)	3.8
FSC threshold	0.143
Map sharpening <i>B</i> factor (Å <sup>2</sup> )	
Model composition	
Non-hydrogen atoms	16305
Protein residues	1958
Ligands	BMA:6 NAG:73 MAN:6
<i>B</i> factors (Å <sup>2</sup> )	
Protein	100.6
R.m.s. deviations	
Bond lengths (Å)	0
Bond angles (°)	0.75
Validation	
MolProbity score	2.10
Clashscore	9.0
Poor rotamers (%)	2.3
Ramachandran plot	
Favored (%)	95.5
Allowed (%)	4.5
Disallowed (%)	0.0

---