

Convergent selection in antibody repertoires is revealed by deep learning

Simon Friedensohn, Daniel Neumeier, Tarik A Khan, Lucia Csepregi, Cristina Parola, Arthur R Gorter de Vries, Lena Erlach, Derek M Mason and Sai T Reddy*

Department of Biosystems Science and Engineering, ETH Zürich, 4058 Basel, Switzerland

*To whom correspondence should be addressed. Tel: +41 61 387 33 68; Email: sai.reddy@ethz.ch

SUMMARY

Adaptive immunity is driven by the ability of lymphocytes to undergo V(D)J recombination and generate a highly diverse set of immune receptors (B cell receptors/secreted antibodies and T cell receptors) and their subsequent clonal selection and expansion upon molecular recognition of foreign antigens. These principles lead to remarkable, unique and dynamic immune receptor repertoires¹. Deep sequencing provides increasing evidence for the presence of commonly shared (convergent) receptors across individual organisms within one species²⁻⁴. Convergent selection of specific receptors towards various antigens offers one explanation for these findings. For example, single cases of convergence have been reported in antibody repertoires of viral infection or allergy⁵⁻⁸. Recent studies demonstrate that convergent selection of sequence motifs within T cell receptor (TCR) repertoires can be identified on an even wider scale^{9,10}. Here we report that there is extensive convergent selection in antibody repertoires of mice for a range of protein antigens and immunization conditions. We employed a deep learning approach utilizing variational autoencoders (VAEs) to model the underlying process of B cell receptor (BCR) recombination and assume that the data generation follows a Gaussian mixture model (GMM) in latent space. This provides both a latent embedding and cluster labels that group similar sequences, thus enabling the discovery of a multitude of convergent, antigen-associated sequence patterns. Using a linear, one-versus-all support vector machine (SVM), we confirm that the identified sequence patterns are predictive of antigenic exposure and outperform predictions based on the occurrence of public clones. Recombinant expression of both natural and *in silico*-generated antibodies possessing convergent patterns confirms their binding specificity to target antigens. Our work highlights to which extent convergence in antibody repertoires can occur and shows how deep learning can be applied for immunodiagnostics and antibody discovery and engineering.

1 RESULTS

2 Targeted deep sequencing of the rearranged BCR locus reveals the antibody repertoire of B cells in a given
3 tissue or cell population¹¹. Here we used deep sequencing to analyse the antibody repertoires in the bone
4 marrow of 45 BALB/c mice, which were divided into four separate cohorts immunized with protein antigens
5 of either ovalbumin (OVA), hen egg lysozyme (HEL), blue carrier protein (BCP) or respiratory syncytial virus
6 fusion protein (RSV-F). OVA, HEL and BCP cohorts were further subdivided into groups receiving zero, one,
7 two or three booster immunizations; serum enzyme-linked immunoabsorbance assays (ELISA) confirmed
8 antigen-specific antibody titers in all mice, with mice receiving only a primary immunization exhibiting
9 substantially weaker titers (**Supplementary Table 1**). RNA was extracted in bulk from bone marrow cells and
10 variable heavy (V_H) chain IgG sequencing libraries were prepared using a two-step RT-PCR protocol, which
11 included molecular barcoding for error and bias correction¹². Libraries were sequenced on an Illumina MiSeq,
12 quality processed and aligned, yielding across all mice a total of 243'374 unique combinations of all three V_H
13 complementarity-determining regions (CDRHs) (**Supplementary Fig. 1a, b**). Similar to previous studies³, we
14 observed the occurrence of public clones (defined here as identical CDRH1-CDRH2-CDRH3 amino acid (a.a.)
15 sequences that occur in more than one mouse). However, the occurrence of public clones was not limited by
16 antigen exposure, as on average 16% of the sequences found in a given repertoire were shared with
17 repertoires outside the respective antigen cohort, whereas only ~4% were limited to a specific antigen cohort
18 (**Supplementary Fig. 1c**).

19
20 To evaluate to which extent convergent selection occurs that is beyond the occurrence of public clones, we
21 developed a deep learning workflow which utilizes CDRH1, CDRH2 and CDRH3 and their appropriate
22 sequence combinations as input to a VAE model. Deep generative models using variational inference have
23 recently been used for TCR repertoire modelling and protein fitness prediction¹³⁻¹⁵. Unlike this previous work
24 however, we assume that sequences in latent space are generated by a GMM, which enables the VAE to
25 estimate the densities of different clonal lineages and signatures more accurately (**Fig. 1**). The VAE consists
26 of deep neural networks which are used to encode and decode sequences and are optimized with respect
27 to the GMM prior and their ability to reconstruct the input sequences (**Extended Data Fig. 1**), with similar
28 sequences falling into the same cluster and closely related clusters occupying similar regions in latent space
29 (**see Methods**). Increasing the dimensionality of the latent encoding and the number of clusters improved the
30 reconstruction ability of the model; by using a ten-dimensional latent space with 2,000 clusters, we were able
31 to achieve a reconstruction accuracy of 93.4% (**Extended Data Fig. 2**). We used principal component
32 analysis (PCA) to visualize the latent encodings and found that highly similar sequences did indeed map to
33 the same cluster and areas within the latent space (**Fig. 2a**). The VAE model also captured important
34 information such as CDRH3 length and variable germline segment (V-gene), whereas similar V-gene families
35 were mapped into similar areas of the latent space (**Extended Data Fig. 3**). While visual inspection revealed
36 areas of possible antigen-associated sequence convergence (**Fig. 2a, Extended Data Fig. 4**), we aimed to
37 quantify the amount of convergence by identifying latent clusters that were significantly associated for each
38 respective antigen, and whether these convergent sequences could be used to predict the antigen exposure
39 of a given mouse. Sequences were grouped into their respective clusters and the recoded repertoires were
40 used to train and cross-validate a linear, one-versus-all SVM classifier of antigen exposure (**Supplementary**

41 **Fig. 2).** It is important to note that the VAE allows cluster assignments of unseen data and thus an independent
42 embedding could be trained for each fold. In order to establish a baseline for this workflow, we also trained
43 a linear SVM on the occurrence of public clones (exact CDRH1-CDRH2-CDRH3 a.a. sequence matches),
44 which yielded an accuracy of 42% for prediction of antigen exposure (5-fold cross-validation). In contrast,
45 when using VAE-based cluster assignments and subsequently encoding repertoires based on cluster
46 enrichment, the resulting classifiers were able to achieve a prediction accuracy of over 80% (**Fig. 2b**,
47 **Extended Data Fig. 5 and 6**). We then performed a simple, non-parametric permutation test for each cluster
48 and each cohort at a significance level of $\alpha=0.05$. Bonferroni correction was conducted in order to account
49 for multiple testing, yielding 60 (BCP, 6'664 sequences), 61 (RSV-F, 7'064 sequences), 68 (OVA, 7'389
50 sequences) and 73 (HEL, 9'628 sequences) significantly enriched convergent antibody clusters in each
51 cohort. While the exact number of convergent clusters and sequences slightly varies with the number of latent
52 space clusters and the initialization of the VAE, the anti-correlation between protein antigen size and
53 complexity and identified sequences suggests that convergence may be driven by the presence of a few
54 immunodominant epitopes. Closer inspection revealed that not every mouse expressed all of their respective
55 convergent clusters, but rather a smaller, yet still predictive subset (**Fig. 2c**). Furthermore, mice that only
56 received a primary immunization without any subsequent booster immunization also exhibited a decreased
57 enrichment of convergent clones (**Fig. 2c**, area between dashed red lines), a finding that correlates well with
58 the measured serum titers (**Supplementary Table 1**). Example logos generated by sequences mapping into
59 the same cluster visualize how the VAE model is able to identify diverse and biologically meaningful groupings
60 (**Fig. 2d**). Furthermore, comparing aggregated sequence logos to those generated from single repertoires
61 shows the potential diversity of the convergent sequence space and highlights that convergence is not limited
62 only to highly similar, public CDRHs (**Extended Data Fig. 7**). As an additional, yet simple example of
63 convergence, we observed a frequently occurring asparagine (N) residue in the first position of CDRH3 of
64 multiple RSV-F-associated clusters (**Fig. 2d**).

65

66 In order to confirm that antigen-associated sequence convergence corresponds to antigen binding, we
67 utilised a mammalian cell (hybridoma) antibody surface display and secretion system coupled to CRISPR-
68 Cas9-mediated library integration¹⁶. An antibody library was generated using a small subset of convergent V_H
69 sequences derived from different clusters combined with a variable light chain (V_L) library cloned from cohort-
70 matched mouse repertoires (**Extended Data Fig. 8**). ELISAs performed on the supernatant of the library cell
71 lines demonstrated that these convergent pools possessed cognate antigen specificity (**Fig. 3a**). We then
72 proceeded to more closely investigate V_H variants from the OVA and RSV-F pools through single clone
73 isolation by fluorescence-activated cell sorting (FACS) (**Supplementary Fig. 3**). The antigen-specific binding
74 of monoclonal cell lines was confirmed by flow cytometry (**Fig. 3b**) and ELISA (**Supplementary Fig. 4**) and
75 their V_H chains were identified by sequencing (**Supplementary Fig. 5**). This procedure allowed us to confirm
76 antigen specificity of 6 (out of 6 selected) OVA and 3 (out of 4 selected) RSV-F convergent V_H sequences
77 (**Extended Data Table 1**). V_H chains were able to pair with V_L chains from a different mouse repertoire,
78 additionally highlighting convergence with respect to V_H chain-dominated binding (**Supplementary Tables**
79 **2-5**). While all antigens were associated with a variety of V-gene germlines, we noticed that convergent

80 antibodies were utilizing different V-gene segments in an antigen-dependent manner, highlighting that the
81 original V-gene germline contributes to convergent selection (**Fig. 3f, Extended Data Fig. 9**).

82

83 In order to confirm that antibody sequence variants mapping to the same convergent cluster were also
84 antigen-specific, we recombinantly expressed 12 convergent V_H variants (derived from other mice immunized
85 with the same antigen) from the cluster mapping to one of the confirmed RSV-F binders (RSV3,
86 **Supplementary Fig. 6**). These 12 convergent variants were expressed with the same V_L of RSV3. Flow
87 cytometry confirmed that all 12 of these convergent variants were indeed antigen-specific (**Fig. 3c**). Using
88 standard clonotype definitions of 100% or 80% V_H CDRH3 a.a. identity^{2,4}, only zero or five of the 12 variants,
89 respectively, would have been identified as convergent across repertoires (**Fig. 3d**). In contrast, the VAE
90 model was able to discover variants of RSV3 with as low as 64% CDRH3 a.a. identity (4 out of 11
91 mismatches), verifying the large potential diversity revealed by the previous logo plots (**Fig. 2d, Fig. 3f**).
92 Besides their sequence diversity, these clones also confirmed the large abundance range with confirmed
93 binders being of high, medium and low frequencies in their respective mouse repertoires (**Fig. 3e**).

94

95 Finally, we aimed to understand how well the VAE model is able to generalise to unseen data. To start, we
96 experimentally produced an antibody CDRH3 library of the RSV3 clone through CRISPR-Cas9 homology-
97 directed mutagenesis¹⁷; while the diversity of the library was designed to model decoder-generated
98 sequences of the RSV3 cluster, it also contained fully randomized positions (**Supplementary Fig. 7a**).
99 Screening of the hybridoma antibody library by FACS followed by deep sequencing yielded 19'270 surface-
100 expressed variants of whom 7'470 were confirmed to be antigen-binding (**Supp. Fig. 7b**). When assessing
101 the probabilities of these novel variants under the VAE model, we found that binding CDRH3s were
102 significantly more likely to be generated than non-binding variants (**Extended Data Fig. 10**). However, since
103 the library also contained a.a. that were not observed in nature, most of its variants were less likely to be
104 generated by our model than naturally occurring sequences (**Extended Data Fig. 10, Supplementary. Figure**
105 **8**). Yet, the overlap between the distributions indicated that the VAE model should have been able to generate
106 some of these variants *in silico*. We confirmed this fact by sampling one million latent encodings directly from
107 the respective RSV3 containing cluster of the GMM model. The trained decoder was then used to transform
108 the sampled encodings into distinct position probability matrices from which in turn actual sequences were
109 generated (**Fig. 4a**). This procedure exhaustively sampled the CDRH3 sequence space of RSV3 and yielded
110 5'005 novel, high quality *in silico* variants that did not occur in the original biological training set. Of these
111 variants, 71 were confirmed by the previous library screen to bind RSV-F, while 25 *in silico* variants were
112 found in the non-binding fraction; resulting in an overall binding accuracy of 74%. Again, the non-binding
113 variants were sampled at a much lower rate than binding variants, indicating that the bulk of the *in silico*
114 generated sequences are likely to be antigen-specific as well (**Fig. 4b, Extended Data Table 2**).

115

116 DISCUSSION

117 In summary, we show that wide-scale convergence across a range of antigens occurs in the antibody
118 repertoires of mice. Unlike current approaches used to identify clonotypes¹⁸⁻²⁰, our VAE approach learns the
119 clustering thresholds based on densities of individual sequence motifs found in the data, thereby forming

120 clusters of varying degrees of similarity. Furthermore, our trained encoding neural network is able to identify
121 these motifs in unseen repertoires in a sensible manner, thereby extending currently existing frameworks for
122 immunodiagnostics^{21,22}. Commonly applied methods to detect convergence such as clonotyping based on
123 exact V-gene and J-gene and a CDRH3 similarity threshold (e.g. 80%) are only partly able to recover
124 convergent patterns; in contrast deep learning revealed that convergent antibody sequences utilized various
125 V-genes and J-genes and have dissimilarities of up to >40% in CDRH3. It is therefore likely that the current
126 extent of antigen driven convergence in immune receptor repertoires has been underreported. We are able
127 to discover convergent motifs from sequences buried deep in the repertoire, highlighting again the possibility
128 that as the amount of available sequencing data increases, similar phenomena might be more commonly
129 observed in humans as well. While we focused specifically on the V_H region (previous studies have shown
130 them to be sufficient to determine most clonal relationships²³), we also observed evidence for potential
131 convergence on the V_L region; future work using single-cell sequencing²⁴ may reveal convergent patterns
132 across the paired V_H:V_L binding domain. We show that our deep generative model allows us to exhaustively
133 sample from the naturally occurring sequence space of antibody repertoires, resulting in the *in silico*
134 generation of antibody variants that retain antigen-binding, a procedure that could be used for engineering
135 antibodies with desired therapeutic development properties²⁵. Detection of convergent patterns by deep
136 learning may also enable the discovery of functional and protective antibodies in patients with unique
137 immunological phenotypes (e.g., elite neutralizers of HIV), which could be exploited as immunodiagnostics,
138 therapeutic antibodies or for vaccine immunogen design^{21,26-28}.

139

140 **METHODS**

141 **Mouse immunizations and RNA isolation from bone marrow**

142 Female BALB/c mice (Charles Rivers) of 6-8 weeks old were separated into cohorts (10-12 mice) based on
143 antigen: hen egg lysozyme (HEL, Sigma Aldrich, 62971-F), ovalbumin (OVA, Hyglos, 321001), blue carrier
144 protein (BCP, Pierce, 771300) and respiratory syncytial virus fusion glycoprotein (RSV-F, expressed
145 internally). For HEL, OVA and BCP, mice were immunized with an initial subcutaneous injection of 200 µg
146 antigen (and 20 µg monophosphoryl lipid A (MPLA, Invivogen, Tlrl-mpla) adjuvant. These mice received zero,
147 one, two or three booster injections for which final immunizations (boost 1, 2 or 3) were done with 50 µg
148 antigen (intraperitoneal injection without any adjuvants). Middle immunizations (boost 1 and/or 2) were done
149 with 50 µg antigen and 20 µg MPLA. For RSV-F, all mice were immunized with 2 booster injections, with each
150 of the three injections consisting of 10 µg for RSV-F and 1% Alum (Thermo Scientific, 77161) adjuvant. For
151 all antigens, sequential injections were interspaced by three weeks. All adjuvants and antigens were prepared
152 and aliquoted before the experiments and mixed in a total volume of 150 µL (for RSV-F: 100µl) on the days
153 of the corresponding injection. Mice were sacrificed 10 days (for RSV-F: 14 days) after their final immunization
154 and bone marrow was extracted from femurs of hind legs using 2 mM PBS buffer. The isolated bone marrow
155 was then centrifuged at 400 x g at 4 °C for 5 minutes. The supernatant was removed and 1.25 mL of Trizol
156 (Invitrogen, 15596026) was added. The bone marrow was then homogenized using a 18G*2'' needle (1.2*50
157 mm). 1 mL of the resulting Trizol solution was then frozen at -80 °C until processing. Mouse cohorts and
158 immunization groups are described in **Supplementary Table 1**. RNA extraction was performed as previously
159 described¹². Briefly, 1 mL of the homogenate was used as input for the PureLink RNA Mini Kit (Life

160 Technologies, 12183018A). RNA extraction was then conducted according to the manufacturer's guidelines.

161

162 **Antibody repertoire library preparation and deep sequencing**

163 Antibody variable heavy chain (V_H) libraries for deep sequencing were prepared using a previously established
164 protocol of molecular amplification fingerprinting (MAF), which enables comprehensive error and bias
165 correction¹². Briefly, a first step of reverse transcription was performed on total RNA using a gene-specific
166 primer corresponding to constant heavy region 1 (CH1) of IgG subtypes and with an overhang region
167 containing a reverse unique molecular identifier (RID). Next, multiplex-PCR is performed on first-strand cDNA
168 using a forward primer set that anneals to framework 1 (FR1) regions of V_H and has an overhang region of
169 forward unique molecular identifier (FID) and partial Illumina adapter; reverse primer also contains a partial
170 Illumina sequencing adapter. A final singleplex-PCR step is performed to complete the addition of full Illumina
171 adapters. After library preparation, overall library quality and concentration was determined on a Fragment
172 Analyzer (Agilent). Libraries were then pooled and sequenced on an Illumina MiSeq using the reagent v3 kit
173 (2x300 bp) with 10% PhiX DNA added for quality purposes.

174

175 **Data pre-processing and sequence alignment**

176 Before alignment, the raw FASTQ files were processed by a custom CLC Genomics Workbench 10 script.
177 Firstly, low quality nucleotides were removed using the quality trimming option with a quality limit of 0.05.
178 Afterwards, forward and reverse read pairs were merged and resulting amplicons between 350 and 600 base
179 pairs were kept for further analysis. Pre-processed sequences were then error-corrected and aligned using
180 the previously established MAF bioinformatic pipeline¹². Downstream analysis was then carried out using
181 Python utilizing both the Tensorflow and scikit-learn libraries.

182

183 **Variational autoencoder models of antibody repertoires**

184 Following error and bias correction and alignment of antibody repertoire sequencing data, all discovered
185 combinations of CDRH1, CDRH2 and CDRH3 for each dataset were extracted. In order to process CDRHs
186 of various lengths, sequences were padded with dashes until a certain fixed length (maximum length for each
187 CDRH in the data) was reached. Padded sequences were one-hot encoded, concatenated and used as input
188 into the variational autoencoder (VAE). The VAE model jointly optimizes the ability to reconstruct its input
189 together with a Gaussian mixture model (GMM)-based clustering of the latent space (Fig. 1) according to the
190 following formula:

$$191 \quad \mathcal{L}_{ELBO}(x) = \mathbb{E}_{q(y, z|x)}[\ln p(x, y, z) - \ln q(y, z|x)]$$

192 With:

$$193 \quad p(x, z, y) = p(c)p(z|y)p(x|z)$$

$$194 \quad p(y) \sim \text{Cat}(\pi)$$

$$195 \quad p(z|y) \sim \mathcal{N}(\mu_y, \sigma_y^2 \mathbb{I}_D)$$

196

197 And the following variational approximation of the posterior, where $q(z|x, y)$ is assumed to be distributed
198 according to a gaussian distribution:

199

$$q(y, z|x) = q(y|x)q(z|x, y)$$

200
201
202 Unlike comparable models²⁹, we do not make a mean field approximation when modelling the posterior,
203 which we found to increase model stability. This choice however comes at the expense of considerable
204 increases in computation time, as has been reported before³⁰. Specifically, we encode and decode every
205 input sequence as if it would belong to every cluster (indicated through a one-hot encoded cluster label) using
206 shared weights in every layer. The final contributions to the overall loss are then weighted by the separately
207 predicted probabilities $q(y|x)$, which describe the probability of a sequence belonging to a specific cluster
208 (**Extended Data Fig. 1**). The decoder maps input sequences and concatenated class label into a lower
209 dimensional ($d=10$) space using two dense layer with rectified linear unit (ReLU) activation followed by the
210 final 10-dimensional layer. Sequences are sampled and recreated from the latent space using the decoder.
211 The decoding network (**Extended Data Fig. 1**) employs two separate dense layers with ReLU activation
212 followed by a dense layer with a linear activation, whose output is reshaped and normalized with a softmax
213 activation in order to reconstruct the probabilities of the initial, one-hot encoded CDRHs. The standard
214 categorical cross-entropy loss is used as the reconstruction term. Every VAE model was trained on a single
215 GPU node of the ETH Zurich parallel computing cluster (*Leonhard*). Training consisted of 200 epochs for all
216 models using Adam as the optimization algorithm³¹.

217 218 **Predicting antigen exposure of single antibody repertoires**

219 Repertoire datasets were split into five folds with each fold being approximately balanced in the number of
220 distinct antigen groups and each dataset appearing only once across all folds. This split was then used to
221 perform a cross-validation procedure in which each of the five folds were set aside as a test set once and the
222 remaining four folds were used as training data. For each of the five training/test splits a separate VAE model
223 was learned by combining all sequences across all repertoires from the training set as input. Clustering
224 assignments or sequences from both the training and the test set were then calculated for the trained model.
225 Based on these cluster labels each repertoire was recoded as an n -dimensional vector, where n is the number
226 of possible clusters and the i -th element indicates the number of sequences mapping to the i -th cluster in the
227 given repertoire. These vectors were then used to train and validate linear support vector machines (SVM) in
228 a one-versus-all setting. In order to prevent a more resource-intensive nested cross-validation procedure we
229 decided to not optimize the hyperparameters of the SVMs and instead chose to use the standard values
230 given by scikit-learn's 'SVC' implementation (using a linear kernel). For visualization purposes the results of
231 each cross-validation step were grouped together in one single confusion-matrix (**Fig. 2b**).

232 233 **Identification of convergent antigen-associated sequence clusters**

234 In order to identify convergent antigen-associated sequence clusters from antibody repertoires, we
235 performed a non-parametric permutation test in order to determine whether sequence reads were specifically
236 enriched in one cluster given a specific cohort (**Fig. 2c**). In order to account for multiple testing, Bonferroni
237 correction was applied to all p-values in each cohort. We proceeded by randomly choosing one CDRH1-
238 CDRH2-CDRH3 combination and its cognate full-length variable region from each cluster for further
239 validation. Permutation tests were carried out in Python using *mlxtend* and 1000 permutations each.

240 **Generation of cluster-specific, in-silico variants**

241 Cluster-specific, novel variants were generated in-silico by sampling data points in latent space from a
242 multivariate gaussian distribution, where parameters were given by the respective cluster parameters from
243 the final VAE model. These sampled data points were then fed into the decoding network resulting in position
244 probability matrices for each CDRH (**Fig. 4a**). For each data point a given CDRH1, CDRH2 and CDRH3 was
245 generated. This process was repeated for a million iterations. The log probability of single sequences was
246 approximated by taking the average of 500 samples of the evidence lower bound (ELBO).

247

248 **Hybridoma cell culture conditions**

249 All hybridoma cell lines and libraries were cultivated in high-glucose Dulbecco's Modified Eagle Medium
250 (DMEM; Thermo, 61965026) supplemented with 10% (v/v) heat inactivated fetal bovine serum (FBS; Thermo,
251 26140079), 100 U/ml penicillin/streptomycin (Pen Strep; Thermo, 15140122), 10 mM HEPES buffer (Thermo,
252 15630080) and 50 μ M 2-Mercaptoethanol (Thermo, 21985023). All hybridoma cultures were maintained in
253 cell culture incubators at a constant temperature of 37C in humidified air with 5% CO₂. Hybridoma cells were
254 typically cultured in 10 ml of medium in T-25 flasks (TPP, 90026) and passaged every 48/72h. All hybridoma
255 cell lines were confirmed annually to be *Mycoplasma*-free (Universal Mycoplasma Detection Kit, ATCC, 30-
256 1012K). The hybridoma cell line PnP-mRuby/Cas9 was previously published¹⁷.

257

258 **Generation of antibody libraries by CRISPR-Cas9 homology-directed repair**

259 Candidate V_H genes were ordered from Twist Bioscience as gene fragments, which were resuspended in 25
260 μ l Tris-EDTA, pH 7.4 (Sigma) prior to use. All oligonucleotides as well as crRNAs and tracrRNA used in this
261 study were purchased from Integrated DNA Technologies (IDT) and adjusted to 100 μ M (oligonucleotides)
262 with Tris-EDTA or to 200 μ M (crRNA/tracrRNAs) with nuclease-free duplex buffer (IDT, 11-01-03-01) prior to
263 use. The homology-directed repair (HDR) donor template used throughout this study was based on the
264 pUC57(Kan)-HEL23-HDR homology donor plasmid, previously described^{16,32}. Importantly, two consecutive
265 stop codons were incorporated into the beginning of the coding regions for the V_H and the variable light chain
266 (V_L) sequences in order to avoid library cloning artefacts and background antibody expression due to
267 unmodified parental vector DNA.

268

269 For each candidate V_H, separate HDR-donor V_L libraries were assembled in a stepwise manner by Gibson
270 cloning using the Gibson Assembly Master Mix (NEB, E2611)³³. When necessary, fragments were amplified
271 using the KAPA Hifi HotStart Ready Mix (KAPA Biosystems, K2602) following manufacturer instructions. First,
272 V_H genes were amplified from gene fragments and cloned into the PCR-linearized parental HDR-donor vector
273 (step 1, Extended Data Figure 8). Next, with total bone-marrow RNA of a mouse that was immunized with
274 one of the four respective antigens, reverse transcription was performed using the Maxima Reverse
275 Transcriptase (Thermo, EP0741) with a primer specific for V_L constant region. The resulting cDNA was used
276 to amplify the respective V_L repertoires in PCR reactions using a degenerate multiplex primer mix, previously
277 described³⁴ (**Supplementary Table 6**). V_L repertoires were cloned into the PCR-linearized HDR-donor vector
278 created in step 1 for each candidate V_H library (step 2) and final libraries were assessed in terms of diversity

279 and background clones. Typically, fixed V_H HDR-donor V_L libraries had sizes ranging from 30'000 – 80'000
280 transformants per library.

281

282 PnP-mRuby/Cas9 cells were electroporated with the 4D-Nucleofector System (Lonza) using the SF Cell Line
283 4D-Nucleofector Kit L (Lonza, V4XC-2012) with the program CQ-104. For each HDR-donor library, 10^6 cells
284 were harvested by centrifugation at 125 x g for 10 min, washed with 1 ml of Opti-MEM Reduced Serum
285 Medium (Thermo, 31985-062) and centrifuged again using the same parameters. The cells were finally
286 resuspended in 100 μ l of nucleofection mix containing 500 pmol of crRNA-J/tracrRNA complex and 20 μ g of
287 HDR-donor plasmid (5.9 kb) diluted in SF buffer. Following electroporation, cells were cultured in 1mL of
288 growth media in 24-well plates (Thermo) for two days and moved to 6-well plates (Costar) containing another
289 2 mL of growth media for one additional day.

290

291 **Screening of hybridoma antibody libraries by flow cytometry**

292 Flow-cytometry analysis and FACS of CRISPR-Cas9 modified hybridomas was performed on a BD LSR
293 Fortessa and BD FACS Aria III (BD Biosciences). Flow cytometry data was analyzed using FlowJo V10 (FlowJo
294 LLC). Three days post-transfection, hybridoma cell libraries specific for one antigen were pooled and enriched
295 for antibody-expressing and antigen-specific cells in consecutive rounds of FACS. Typically, the number of
296 sorted cells from the previous enrichment-step was over-sampled by a factor of 40 in terms of the number of
297 labelled cells for the subsequent sorting-step. For labelling, cells were washed with PBS (Thermo, 10010023),
298 incubated with the labelling antibodies or antigen for 30 min on ice protected from light, washed two times
299 with PBS and analyzed or sorted. The labelling reagents and working concentrations are listed in
300 **Supplementary Table 7**. For cell numbers different from 10^6 , the amount of antibody/antigen as well as the
301 incubation volume were adjusted proportionally. For labelling of RSV-F-specific cells, a two-step labelling
302 procedure was necessary due to the indirect labelling of cells with RSV-F-biotin/Streptavidin-AlexaFluor647.

303

304 **Genotyping of monoclonal hybridoma cell lines**

305 Genomic DNA of single-cell sorted and expanded hybridoma clones was isolated from 5×10^5 cells, which
306 were washed with PBS and resuspended in QuickExtract DNA Extraction Solution (Epicentre, QE09050).
307 Cells were incubated at 68 °C for 15 min and 95 °C for 8 min and the integrated synthetic V_L - C_k -2A- V_H
308 antibody region was PCR-amplified with flanking primers 5'-CATGTGCCTTTTCAGTGCTTTCTC-3' and 5'-
309 CTAGATGCCTTTCTCCCTTGACTC-3' that were specific for the 5' and 3' homology arms. From this single
310 amplicon, both V_H and V_L regions could be Sanger-sequenced using primers 5'-TGACCTTCTCAAGTTGGC-
311 3' and 5'-GAAAACAACATATGACTCCTGTCTTC-3', respectively (Microsynth).

312

313 **Measurement of antibody specificity by ELISA (cell culture supernatant)**

314 Standard sandwich enzyme-linked immunoabsorbance assays (ELISAs) were performed to measure the
315 specificity of hybridoma cell line supernatants containing secreted IgG. High binding 96-well plates (Costar,
316 CLS3590) were coated over night with 4 μ g/ml of antigen in PBS at 4C. The plates were then blocked for two
317 hours at room temperature with PBS containing 2 % (m/v) non-fat dried milk powder (AppliChem, A0830) and
318 0.05 % (v/v) Tween-20 (AppliChem, A1389). After blocking, plates were washed three times with PBS

319 containing 0.05 % (v/v) Tween-20 (PBST). Cell culture supernatants were 0.2 μ m sterile-filtrated (Sartorius,
320 16534) and serially diluted across the plate (1:3 steps) in PBS supplemented with 2 % (m/v) milk (PBSM),
321 starting with non-diluted supernatants as the highest concentrations. Plates were incubated for one hour at
322 room temperature and washed three times with PBST. HRP-conjugated rat monoclonal [187.1] anti-mouse
323 kappa light chain antibody (abcam ab99617) was used as secondary detection antibody, concentrated at 0.7
324 μ g/ml (1:1500 dilution from stock) in PBSM. Plates were incubated at room temperature for one hour again,
325 followed by three washing steps with PBST. ELISA detection was performed using the 1-Step Ultra TMB-
326 ELISA Substrate Solution (Thermo, 34028) and reaction was terminated with 1 M H₂SO₄. Absorption at 450
327 nm was measured with the Infinite 200 PRO NanoQuant (Tecan) and data were analyzed using Prism V8
328 (Graphpad).

329

330 **Measurement of antibody serum titers by ELISA**

331 Serum titers were measured in a similar manner as described above with the following exceptions: (1) Plates
332 were coated with either 10 μ g/mL (OVA, BCP, HEL) or 2 μ g/mL of antigen (RSV-F) dissolved in PBS. (2) OVA,
333 BCP and HEL serum ELISAs were blocked with 300 μ L/well of PBS with 3% BSA, 3%FBS, 0.05% Tween and
334 0.05% Proclin and were incubated overnight at 4 °C.

335

336 **RSV3 CDRH3 library generation**

337 RSV3 CDRH3 libraries were generated following CRISPR-Cas9 homology-directed mutagenesis, as
338 previously described¹⁷. Briefly, a single-stranded oligonucleotide (ssODN) encoding a nucleotide sequence
339 that put the endogenous CDRH3 out of frame and contained a highly efficient CRISPR targeting sequence
340 was incorporated into the genomic CDRH3 locus of the RSV3 cell line by CRISPR/Cas9-mediated HDR, using
341 reagents crRNA DN_RSV3_H3-3 and 500 pmol of DN_RSV3-OOF ssODN HDR-template (RSV3-OOF cell line,
342 **Supplementary Fig. 7a, Supplementary Table 7**). Next, 4 x 10⁶ RSV3-OOF cells were transfected with
343 crRNA-DM1 and 500pmol of ssODN encoding a CDRH3 library template per 1 x 10⁶ cells to generate the
344 RSV3 CDRH3 in silico library. Transfected cells were subsequently sorted in two consecutive steps for
345 antibody expression and specificity/negativity towards RSV-F (as described above) to enrich for a pure RSV-
346 F-specific cell population and an RSV-F-negative cell population.

347

348 **Deep sequencing of RSV3 CDRH3 libraries**

349 Sample preparation for deep sequencing was performed following a previously established two-step primer
350 extension protocol [23]. Genomic DNA was extracted from 5 x 10⁶ cells of IgG+, IgG+/RSVF+ and IgG+/RSVF-
351 enriched cell populations using the PureLink™ Genomic DNA Mini Kit (Thermo, K182001). Extracted genomic
352 DNA was amplified in a first PCR using a forward primer binding to the beginning of FR3 and a reverse primer
353 binding to the intronic region located ~40 bp downstream of the end of the J-gene.

354 PCRs were performed with Q5 High-Fidelity DNA polymerase (NEB, M0491L) in 8 parallel 50ul reactions with
355 the following cycle conditions: 98 °C for 30 s; 20 cycles of 98 °C for 10 s, 64 °C for 20 s, 72 °C for 20 s; final
356 extension 72 °C for 2 min; 4 °C storage. PCR products were subsequently concentrated using the DNA Clean
357 and Concentrator Kit (Zymo, D4013) followed by 1.2X SPRIselect (Beckman Coulter, B22318) left-sided size
358 selection. Total PCR1 product (~700 ng) was amplified in a second PCR step, which added extension-specific

359 full-length Illumina adapter sequences to each library by choosing 3 different index reverse primers (DM142-
360 144, using DM125 as the forward primer). Cycle conditions were as follows: 98 °C for 30 s; 2 cycles of 98 °C
361 for 10 s, 40 °C for 20 s, 72 °C for 1 min; 15 cycles of 98 °C for 10 s, 65 °C for 20 s, 72 °C for 30 s; final
362 extension 72 °C for 2 min; 4 °C storage. PCR2 products were subsequently concentrated again using the
363 DNA Clean and Concentrator Kit and run on a 1% agarose gel. Product bands of the correct size (~400 bp)
364 were gel-purified using the Zymoclean™ Gel DNA Recovery Kit (Zymo, D4008) and subsequently analyzed
365 by capillary electrophoresis (Fragment analyzer, Agilent). After quantitation, libraries were pooled accordingly
366 and sequenced on a MiSeq System (Illumina) with the paired-end 2x300bp kit.

367

368 **AUTHOR CONTRIBUTIONS**

369 S.F., D.N., T.A.K., D.M. and S.T.R. designed experiments. T.A.K, A.R.G.D.V., L.C. performed mouse-related
370 experiments. S.F., A.R.G.D.V. and L.C. prepared sequencing libraries. D.N., C.P. conducted preliminary
371 experiments. D.N, L.E. conducted experimental antibody library screening. S.F. was responsible for the
372 bioinformatics pipeline. S.F. analyzed deep sequencing data and developed machine learning and deep
373 learning models. S.F., D.N. prepared figures. S.F., D.N., and S.T.R. wrote the paper.

374

375 **ACKNOWLEDGMENTS**

376 We acknowledge the ETH Zurich D-BSSE Single Cell Unit and the Genomics Facility Basel for excellent
377 support and assistance. This work was supported by the European Research Council Starting Grant
378 679403 (to S.T.R.) and ETH Zurich Research Grants (to S.T.R.) and NCCR Molecular Systems Engineering
379 (to S.T.R). The professorship of S.T.R. is supported by an endowment from the S. Leslie Misrock
380 Foundation.

381

382 **COMPETING INTERESTS**

383 ETH Zurich has filed for patent protection on the technology described herein, and S.F. and S.T.R. are
384 named as co-inventors on this patent (United States Patent and Trademark Office Provisional Application:
385 62/843,010).

386

387 **DATA AVAILABILITY**

388 The raw FASTQ files from deep sequencing that support the findings of this study will be deposited
389 (following peer-review and publication) in the Sequence Read Archive (SRA) with the primary accession
390 code(s) <code(s) (<https://www.ncbi.nlm.nih.gov/sra>)>. Additional data that support the findings of this study
391 are available from the corresponding author upon reasonable request.

392

393 **CODE AVAILABILITY**

394 Deep learning models were built in Python v3.6.4 using TensorFlow v1.12.0. Models will be available on a
395 github repository (following peer-review and publication).

396

REFERENCES

- 1 Wang, B. *et al.* Functional interrogation and mining of natively paired human VH:VL antibody repertoires. *Nature biotechnology* **36**, 152-155, doi:10.1038/nbt.4052 (2018).
- 2 Briney, B., Inderbitzin, A., Joyce, C. & Burton, D. R. Commonality despite exceptional diversity in the baseline human antibody repertoire. *Nature* **566**, 393-397, doi:10.1038/s41586-019-0879-y (2019).
- 3 Greiff, V. *et al.* Systems Analysis Reveals High Genetic and Antigen-Driven Predetermination of Antibody Repertoires throughout B Cell Development. *Cell Rep* **19**, 1467-1478, doi:10.1016/j.celrep.2017.04.054 (2017).
- 4 Soto, C. *et al.* High frequency of shared clonotypes in human B cell receptor repertoires. *Nature* **566**, 398-402, doi:10.1038/s41586-019-0934-8 (2019).
- 5 Croote, D., Darmanis, S., Nadeau, K. C. & Quake, S. R. High-affinity allergen-specific human antibodies cloned from single IgE B cell transcriptomes. *Science* **362**, 1306-1309, doi:10.1126/science.aau2599 (2018).
- 6 Setliff, I. *et al.* Multi-Donor Longitudinal Antibody Repertoire Sequencing Reveals the Existence of Public Antibody Clonotypes in HIV-1 Infection. *Cell host & microbe* **23**, 845-854 e846, doi:10.1016/j.chom.2018.05.001 (2018).
- 7 Parameswaran, P. *et al.* Convergent antibody signatures in human dengue. *Cell host & microbe* **13**, 691-700, doi:10.1016/j.chom.2013.05.008 (2013).
- 8 Truck, J. *et al.* Identification of antigen-specific B cell receptor sequences using public repertoire analysis. *Journal of immunology* **194**, 252-261, doi:10.4049/jimmunol.1401405 (2015).
- 9 Glanville, J. *et al.* Identifying specificity groups in the T cell receptor repertoire. *Nature* **547**, 94-98, doi:10.1038/nature22976 (2017).
- 10 Dash, P. *et al.* Quantifiable predictive features define epitope-specific T cell receptor repertoires. *Nature* **547**, 89 (2017).
- 11 Friedensohn, S., Khan, T. A. & Reddy, S. T. Advanced Methodologies in High-Throughput Sequencing of Immune Repertoires. *Trends in Biotechnology* **35**, 203-214, doi:10.1016/j.tibtech.2016.09.010 (2017).
- 12 Khan, T. A. *et al.* Accurate and predictive antibody repertoire profiling by molecular amplification fingerprinting. *Sci Adv* **2**, e1501371, doi:10.1126/sciadv.1501371 (2016).
- 13 Davidsen, K. *et al.* Deep generative models for T cell receptor protein sequences. *Elife* **8**, doi:10.7554/eLife.46935 (2019).
- 14 Riesselman, A. J., Ingraham, J. B. & Marks, D. S. Deep generative models of genetic variation capture the effects of mutations. *Nature methods* **15**, 816-822, doi:10.1038/s41592-018-0138-4 (2018).
- 15 Sinai, S., Kelsic, E., Church, G. M. & Nowak, M. A. Variational auto-encoding of protein sequences. *arXiv preprint arXiv:1712.03346* (2017).
- 16 Pogson, M., Parola, C., Kelton, W. J., Heuberger, P. & Reddy, S. T. Immunogenomic engineering of a plug-and-(dis)play hybridoma platform. *Nature communications* **7**, 12535, doi:10.1038/ncomms12535 (2016).
- 17 Mason, D. M. *et al.* High-throughput antibody engineering in mammalian cells by CRISPR/Cas9-mediated homology-directed mutagenesis. *Nucleic acids research* **46**, 7436-7449, doi:10.1093/nar/gky550 (2018).
- 18 Gupta, N. T. *et al.* Hierarchical Clustering Can Identify B Cell Clones with High Confidence in Ig Repertoire Sequencing Data. *Journal of immunology* **198**, 2489-2499, doi:10.4049/jimmunol.1601850 (2017).
- 19 Glanville, J. *et al.* Naive antibody gene-segment frequencies are heritable and unaltered by chronic lymphocyte ablation. *Proceedings of the National Academy of Sciences of the United States of America* **108**, 20066-20071, doi:10.1073/pnas.1107498108 (2011).
- 20 Greiff, V., Miho, E., Menzel, U. & Reddy, S. T. Bioinformatic and Statistical Analysis of Adaptive Immune Repertoires. *Trends in immunology* **36**, 738-749, doi:10.1016/j.it.2015.09.006 (2015).

- 21 Greiff, V. *et al.* A bioinformatic framework for immune repertoire diversity profiling enables detection
of immunological status. *Genome Med* **7**, 49, doi:10.1186/s13073-015-0169-8 (2015).
- 22 Emerson, R. O. *et al.* Immunosequencing identifies signatures of cytomegalovirus exposure history
and HLA-mediated effects on the T cell repertoire. *Nat Genet* **49**, 659-665, doi:10.1038/ng.3822
(2017).
- 23 Zhou, J. Q. & Kleinstein, S. H. Cutting edge: ig H chains are sufficient to determine most B cell
clonal relationships. *The Journal of Immunology* **203**, 1687-1692 (2019).
- 24 Laserson, U. *et al.* High-resolution antibody dynamics of vaccine-induced immune responses.
Proceedings of the National Academy of Sciences of the United States of America **111**, 4928-4933,
doi:10.1073/pnas.1323862111 (2014).
- 25 Jain, T. *et al.* Biophysical properties of the clinical-stage antibody landscape. *Proceedings of the
National Academy of Sciences of the United States of America* **114**, 944-949,
doi:10.1073/pnas.1616408114 (2017).
- 26 Jardine, J. G. *et al.* HIV-1 broadly neutralizing antibody precursor B cells revealed by germline-
targeting immunogen. *Science* **351**, 1458-1463, doi:10.1126/science.aad9195 (2016).
- 27 Sesterhenn, F. *et al.* Boosting subdominant neutralizing antibody responses with a computationally
designed epitope-focused immunogen. *PLoS Biol* **17**, e3000164, doi:10.1371/journal.pbio.3000164
(2019).
- 28 Horns, F., Dekker, C. L. & Quake, S. R. Memory B Cell Activation, Broad Anti-influenza Antibodies,
and Bystander Activation Revealed by Single-Cell Transcriptomics. *Cell Reports* **30**, 905-913. e906
(2020).
- 29 Jiang, Z., Zheng, Y., Tan, H., Tang, B. & Zhou, H. Variational Deep Embedding: An Unsupervised
and Generative Approach to Clustering.
- 30 Kingma, D. P., Mohamed, S., Rezende, D. J. & Welling, M. in *Advances in neural information
processing systems*. 3581-3589.
- 31 Kingma, D. P. & Ba, J. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*
(2014).
- 32 Parola, C. *et al.* Antibody discovery and engineering by enhanced CRISPR-Cas9 integration of
variable gene cassette libraries in mammalian cells. *MAbs* **11**, 1367-1380,
doi:10.1080/19420862.2019.1662691 (2019).
- 33 Gibson, D. G. *et al.* Enzymatic assembly of DNA molecules up to several hundred kilobases. *Nature
methods* **6**, 343-U341, doi:10.1038/Nmeth.1318 (2009).
- 34 Reddy, S. T. *et al.* Monoclonal antibodies isolated without screening by analyzing the variable-gene
repertoire of plasma cells. *Nature biotechnology* **28**, 965-969, doi:10.1038/nbt.1673 (2010).
- 35 Lefranc, M. P. Nomenclature of the human immunoglobulin heavy (IGH) genes. *Exp Clin
Immunogenet* **18**, 100-116, doi:49189 (2001).

FIGURES

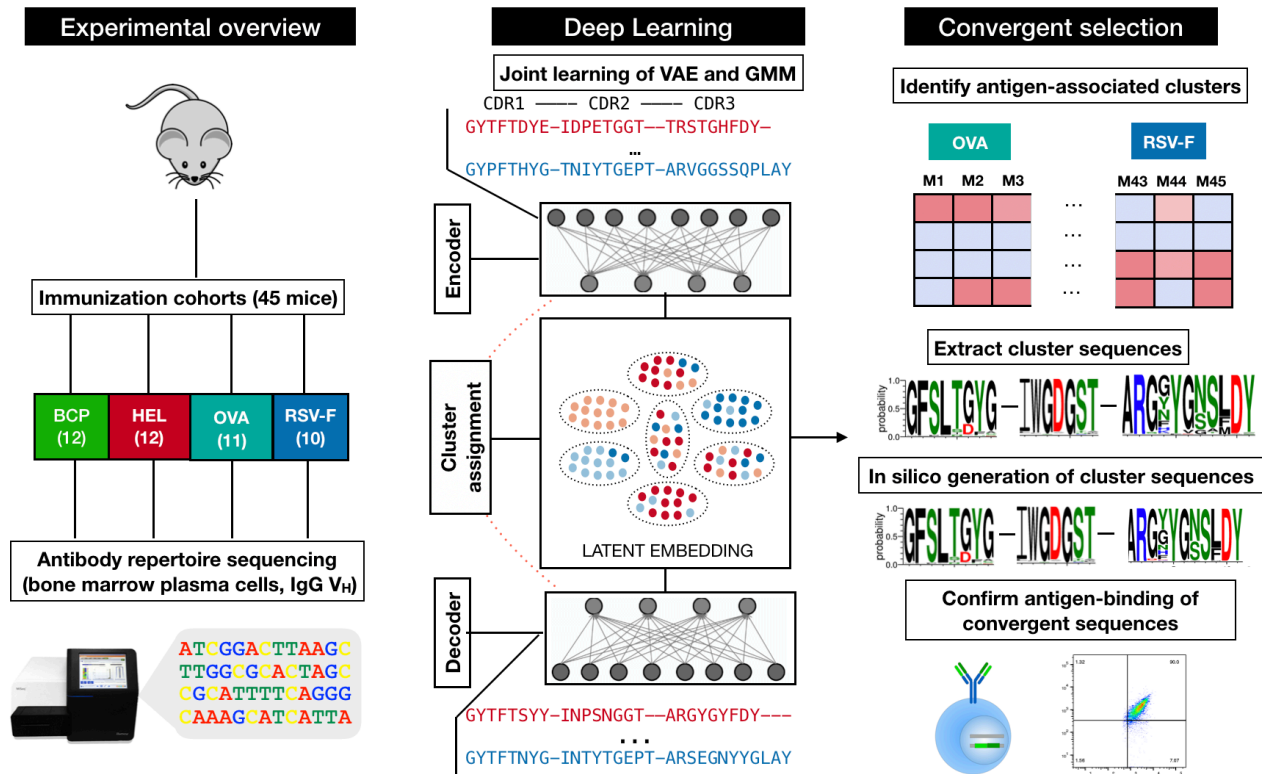


Figure 1. Schematic overview of deep learning on antibody repertoires of immunized mice.

Antibody repertoires from the bone marrow of mice immunized with various antigens are deep sequenced. Antibody sequences are then used to train a variational autoencoder (VAE) with a Gaussian mixture model (GMM) clustering of the latent space. The VAE model is able to both generate novel sequences and assign input sequences to distinct clusters based on their latent space embedding. Cluster assignments are used to identify convergent sequences that are heavily enriched in a specific repertoire or antigen cohort. Natural and in-silico generated sequences from antigen-associated clusters are expressed as full-length IgG and verified as binding antibodies.

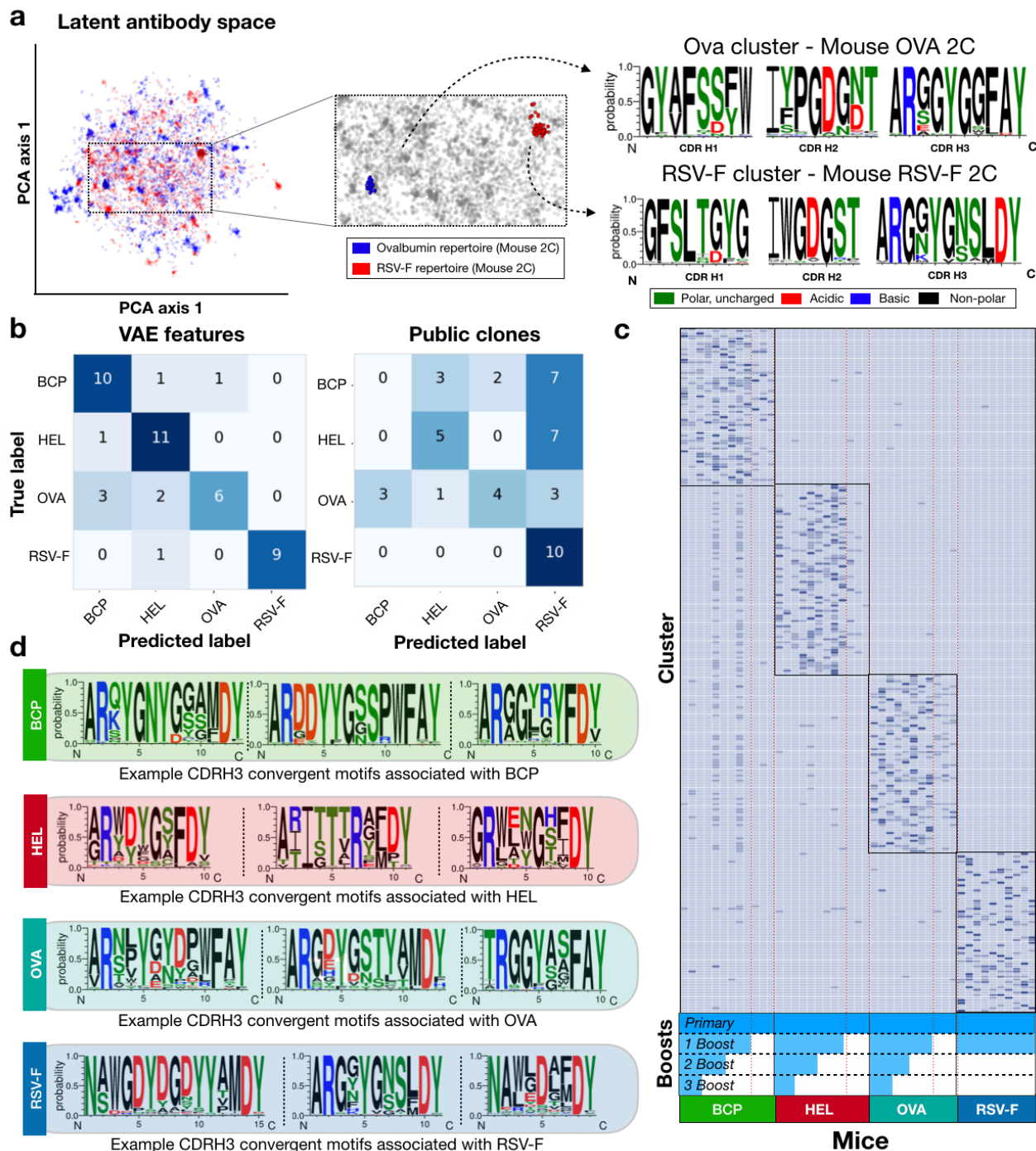


Figure 2. Identification and characterization of convergent antigen-associated sequences.

a, Ten-dimensional latent space of two antibody repertoires visualized by principal component analysis (PCA). Blue and red dots indicate sequences belonging to one OVA (2C) and RSV-F (2C) repertoire, respectively. Enlarged area highlights two learned clusters only containing sequences specific to one repertoire and their respective sequence motifs. b, Antibody repertoires are transformed into vectors based on the learned sequence clusters in latent space. Recoded vectors are used as input for a linear support vector machine (SVM) classifier of antigen exposure. Confusion matrices show the aggregated prediction results of each model during 5-fold cross-validation using the cluster labels and raw sequences as features. c, Heatmap contains all predictive and convergent sequence clusters for each cohort. Dashed red line indicates mice that only received the primary immunization. d, Example sequence logos of convergent clusters found in each antigen cohort.

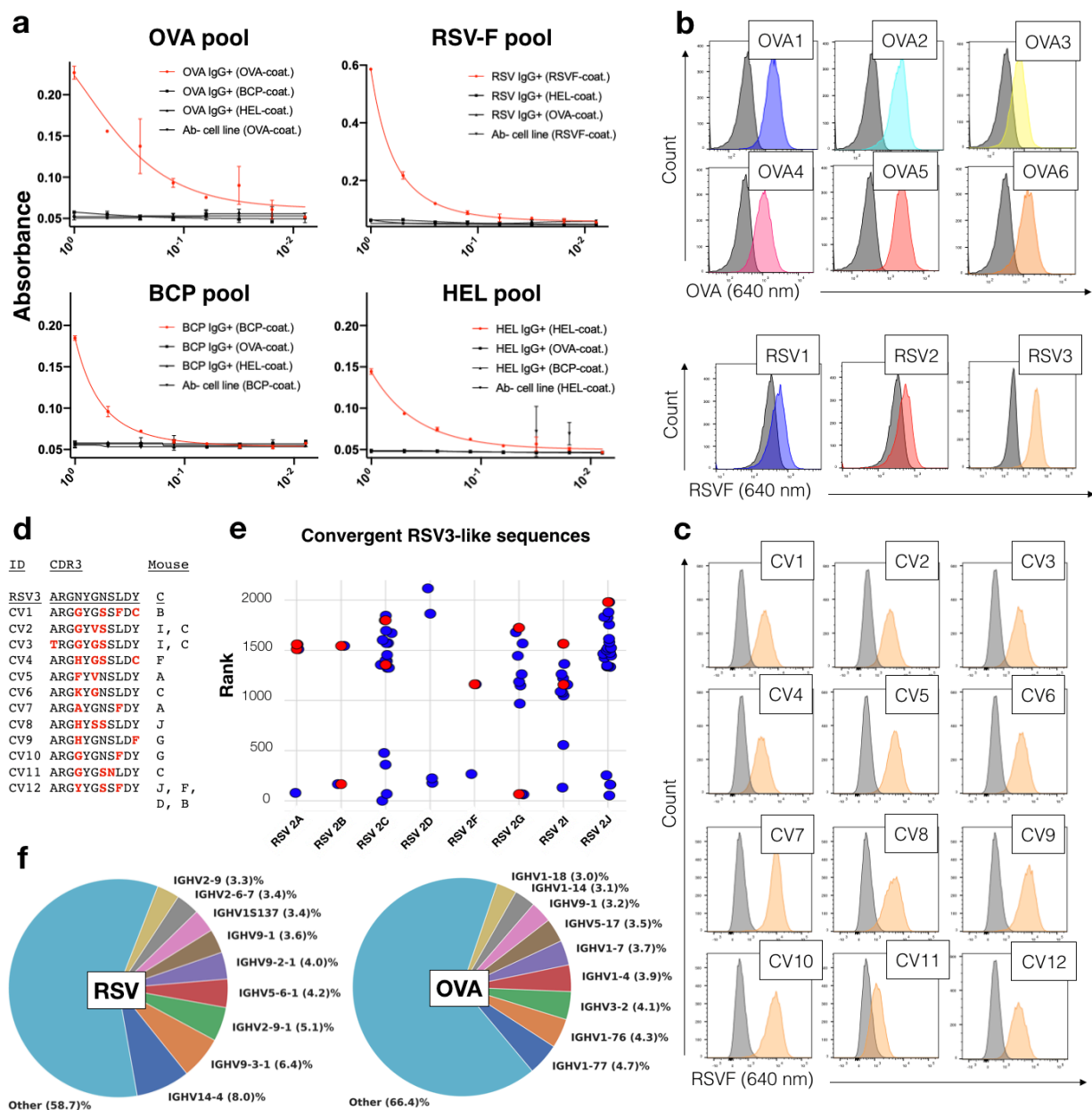


Figure 3. Convergent clusters contain antigen-specific antibodies.

a, Dose-dependent absorbance curves of supernatant prepared from hybridoma cells expressing antibodies with convergent variable heavy (V_H) chain pools for each antigen. b, Flow cytometry histograms of six monoclonal cell populations each utilizing a different convergent OVA-associated or RSV-F associated V_H . Grey histograms represent negative controls, colored histograms show the convergent antibodies. c, Flow cytometry histograms of 12 monoclonal cell populations of convergent variants (CV), which use a different V_H sequence from the same cluster as RSV3. d, Table shows the CDRH3s of the selected CVs and the RSV-F immunized mouse repertoire in which they were found. Red letters indicate differences to the initially discovered sequence RSV3 sequence. e, Scatterplot shows the frequency-rank distributions per mouse repertoire of CVs from RSV3 cluster. Red dots highlight V_H confirmed to be binding in c. e, Pie charts show the nine most utilized V-gene germlines in convergent clones for both RSV-F and OVA.

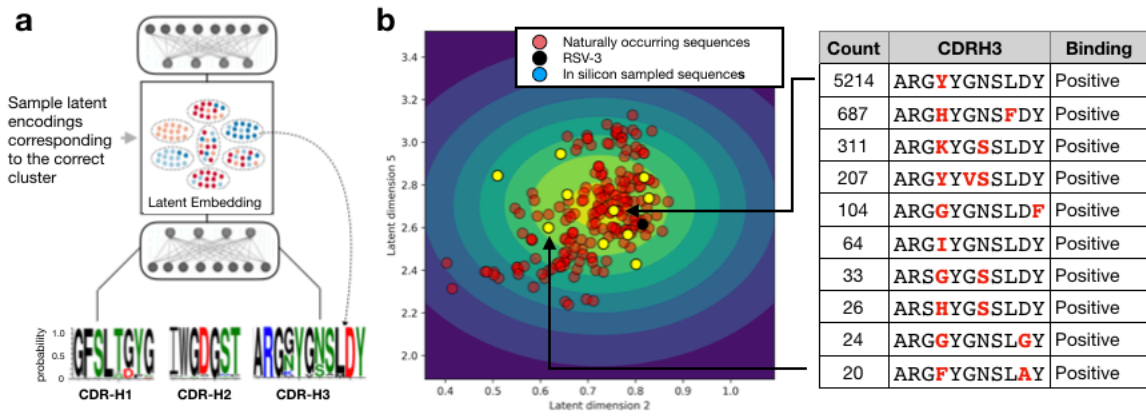


Figure 4. Deep generative modelling and in silico antibody sequence generation.

a, Schematic deep generative modeling of antibody sequence space: a cluster is either chosen or randomly sampled and based on the parameters chosen, a random sample is drawn from a multivariate normal distribution. The encoder then translates the encoding into a multivariate multinomial distribution from which a novel sequence is sampled. b, Scatter plot shows the two latent naturally occurring variants, yellow dots show the ten most frequently in-silico sampled encodings that were confirmed to be binding antibodies. The table on the right shows their CDRH3 sequence and its count after 1'000'000 samples. Red letters indicate differences to the initial biological sequence (RSV3, shown in black).