

Arora et al. (2019) Repertoire-Based Diagnostics Using Statistical Biophysics

Repertoire-Based Diagnostics Using Statistical Biophysics

Rohit Arora, PhD¹, Joseph Kaplinsky, PhD², Anthony Li, MS,¹ and Ramy Arnaout, MD, DPhil^{1,3*}

¹Division of Clinical Pathology, Department of Pathology, Beth Israel Deaconess Medical Center, Boston, MA 02215

²Department of Health Technology, DTU Healthtech, Technical University of Denmark, Produktionstorvet, 2800, Kongens Lyngby, Denmark.

³Harvard Medical School, 25 Shattuck St, Boston, MA 02115

*To whom correspondence should be addressed at rarnaout@gmail.com

Arora et al. (2019) Repertoire-Based Diagnostics Using Statistical Biophysics

9 Abstract

10 A fundamental challenge in immunology is diagnostic classification based on repertoire se-
 11 quence. We used the principle of maximum entropy (MaxEnt) to build compact representations
 12 of antibody (IgH) and T-cell receptor (TCR β) CDR3 repertoires based on the statistical biophysi-
 13 cal patterns latent in the frequency and ordering of repertoires' constituent amino acids. This
 14 approach results in substantial advantages in quality, dimensionality, and training speed com-
 15 pared to MaxEnt models based solely on the standard 20-letter amino-acid alphabet. De-
 16 scriptor-based models learn patterns that pure amino-acid-based models cannot. We demon-
 17 strate the utility of descriptor models by successfully classifying influenza vaccination status
 18 (AUC=0.97, $p=4\times 10^{-3}$), requiring only 31 samples from 14 individuals. Descriptor-based MaxEnt
 19 modeling is a powerful new method for dissecting, encoding, and classifying complex reper-
 20 toires.

21 Introduction

22 A major challenge in systems immunology is determining how to describe the sequence-level
 23 heterogeneity of antibody (immunoglobulin; Ig) and T-cell receptor (TCR) repertoires in ways
 24 that facilitate the identification of meaningful patterns. Sequence-frequency distributions—for
 25 example, counts of unique IgH or TCR β CDR3s—are commonly used but not ideal for interper-
 26 sonal comparisons, since repertoires from different people are largely disjoint (Robins et al.,
 27 2010; Arnaout et al., 2011). Motif-frequency distributions, which count how often each of the 20^n
 28 possible n -mers appears in a repertoire (for some choice of n), are more likely to overlap be-
 29 tween individuals, but may fail to detect probabilistic or higher-order patterns and are subject to
 30 sampling-related bias unless n is small. Comparisons of frequency distributions between reper-
 31 toires from different individuals have yielded important insights (Parameswaran et al., 2013;
 32 Kaplinsky et al., 2014; Emerson et al., 2017; Sun et al., 2017) but the limitations of this ap-
 33 proach suggest a need for complementary methods. One such method is maximum-entropy
 34 (MaxEnt) modeling (Fig. 1).

35 MaxEnt models, which were first developed for statistical physics and information theory
 36 (Jaynes, 1957), can be used to describe repertoires (or other complex ensembles of proteins,
 37 nucleic acids, etc.) in terms of constraints called *biases* that determine the ways in which a giv-
 38 en repertoire differs from a uniform distribution of sequences (Yeo and Burge, 2004; Russ et al.,
 39 2005; Seno et al., 2008; Mora et al., 2010; Marks et al., 2011). Given a set of *features*—for ex-
 40 ample, the frequencies of the 20 amino acids and the $20^2=400$ nearest-neighbor amino-acid
 41 pairs (“neighbors” being defined as contiguous N-to-C-terminus amino acids)—a MaxEnt model
 42 describes the degree to which each feature is biased away from its value in a uniform repertoire,
 43 taking all the other biases into account. For example, the bias for the pair cysteine-alanine (CA)
 44 describes the extent to which the frequency of CA in the repertoire differs from what would be
 45 expected given the frequencies of the individual amino acids C and A, the pairs XC and AX (for
 46 all amino acids X), and so on. MaxEnt models deconvolute the hundreds or thousands of inter-
 47 actions among features into separate components, which then together govern the generation
 48 of the observed sequence- and motif-frequency distributions. Thus MaxEnt models can be
 49 thought of as capturing the underlying generative structure of the repertoire.

50 MaxEnt modeling of IgH and TCR β CDR3s, as well as of other protein families, has shown that
 51 the frequencies of a single set of neighboring amino-acid pairs capture a remarkable amount of
 52 information (Russ et al., 2005; Seno et al., 2008; Mora et al., 2010; Marks et al., 2011), but not

Arora et al. (2019) Repertoire-Based Diagnostics Using Statistical Biophysics

all of it (Bialek and Ranganathan, 2007). Additional sets of pairs—for example, second-, third-, or fourth-nearest neighbors (Mora et al., 2010)—add precision but at the cost of a substantial increase in the number of model parameters (400 per set of pairs). This increase can affect the coverage per feature (the total number of instances of the features in the sample divided by the number of features), model quality and interpretability, and training time. The root of the problem is that the amino-acid alphabet has 20 letters: as a result, parameters, data, and computational requirements scale roughly as powers of 20. The alphabet also causes a second important problem: letters in and of themselves, while a familiar and useful shorthand, lack information about similarities and differences among the multi-faceted biophysical entities they represent—the amino acids—for example, that A is more like glycine (G) than tyrosine (Y)—that may well contain meaningful patterns that are not obvious from, or captured by, the shorthand alone.

These two problems can be addressed simultaneously by swapping the traditional amino-acid alphabet for a smaller set of *descriptors* derived from amino acids' physicochemical properties, especially for pairs and higher-order associations (e.g. consecutive triples) (Fig. 1a). Over two dozen lipophilic (e.g., hydrophobicity), steric (e.g. molecular weight) and electrical (e.g. charge) properties have been precisely measured (Sandberg et al., 1998; Kim et al., 2016). These properties have been shown to correlate with each other, such that the first few principal components (PCs) explain a majority of the overall variance (Hellberg et al., 1987; Sandberg et al., 1998). These PCs are natural candidates for a reduced alphabet: they define orthogonal dimensions of a continuous space in which the discrete amino acids are embedded (Fig. 1a). Whereas in “letter space” there is no concept of distance between amino acids, in “descriptor space” amino acids with similar properties are closer together (e.g., with A nearer G than Y) (Fig. 1b). Such embeddings have been explored in immune-repertoire analysis (Greiff et al., 2017; Ostmeier et al., 2017, 2019) and other contexts (Dosztányi and Torda, 2001; Walter et al., 2005; Susko and Roger, 2007; Stephenson and Freeland, 2013). We investigated whether descriptor-based MaxEnt models of IgH and TCR β CDR3 repertoires could improve on models based on amino acids alone by allowing more data per parameter (less sampling error), shorter training time, and better interpretability (Fig. 1c-d), in principle leading to better models useful for classification of states of health and disease.

82 Results

83 Using 26 measurements carried out on the 20 standard amino acids, we derived five biophysical
 84 descriptors that together explained 92% of the variance in amino acids' physicochemical proper-
 85 ties. Each descriptor is a PC, i.e. a linear combination of the measurements. The first three de-
 86 scriptors corresponded roughly to surface area/chromatographic properties (explaining 41% of
 87 the overall variance), van der Waals volume (25%), and charge (14%) and together explained
 88 79% of variance, an increase over the 68% previously reported for the first three descriptors de-
 89 rived from measures of both the standard and additional non-canonical amino acids (Sandberg
 90 et al., 1998).

91 We trained amino-acid- and descriptor-based MaxEnt models on representative IgH and TCR β
 92 CDR3 repertoires (Fig. 2) and asked which type of model better described test sets of CDR3
 93 sequences set aside from each repertoire, using a nearest- and next-nearest-neighbor amino-
 94 acid model as the benchmark (Methods) (Mora et al., 2010). We compared this benchmark to
 95 two descriptor models: one that fit similar positional information but with fewer parameters, and
 96 one that fit more positional information with a more similar number of parameters. To score
 97 these comparisons, we calculated the (logarithm of the) relative probability that each sequence
 98 σ in the relevant test set belonged to its repertoire according to each of the two models (M_d , de-
 99 scriptor model; M_a , amino-acid model):

$$\ln \frac{p(\sigma|M_d)}{p(\sigma|M_a)}$$

100 and calculated the percent of sequences for which each model was a better fit. A score of 100%
 101 for a given model meant that that model gave a higher probability for every sequence in the test
 102 set. As validation, we confirmed that IgH models scored >99% of IgH sequences better than
 103 TCR β models (Fig. 3a, left), and TCR β models scored >99% of TCR β sequences better than
 104 IgH models (Fig. 3a, right).

105 Test 1: Similar positional information. We first compared models that incorporated similar posi-
 106 tional information: single-amino-acid positions and nearest- and next-nearest neighbor pairs
 107 (see Methods). The amino-acid models required $2 \times 20^2 = 800$ parameters to capture the pairwise
 108 information vs. just $2 \times 5^2 = 50$ parameters for the descriptor models (for each of IgH and TCR β).
 109 We predicted that amino-acid models would outperform descriptor models on this test, since for
 110 every pair of positions the amino-acid model should have a slight edge, given that descriptors

Arora et al. (2019) Repertoire-Based Diagnostics Using Statistical Biophysics

capture only 92% of the variance in amino acids' biophysical properties. Thus this test was expected to provide an estimate of the cost of swapping alphabets. As predicted, amino-acid models outperformed descriptor models, by a wide margin: 94.2% to 5.8% for IgH (Fig. 3b, left) and 99.6% to 0.4% for TCR β (Fig. 3b, right). The median sequence had a probability that was ~240 (IgH) and ~87,000 (TCR β) times as high according to the amino-acid model as according to the descriptor model. For the amino-acid models, sequences from the final samples often contained the canonical CDR3 stems (see Discussion), but these were rare for final samples from these simple descriptor models.

Test 2: Similar numbers of parameters. A primary motivation for developing descriptor models is their ability to capture information at a given set of positions with fewer parameters than amino-acid models; the corollary is that for a given number of parameters, descriptor models can capture more positional information. Specifically, for the 400 parameters amino-acid models require to capture information about nearest-neighbor pairs, descriptor models can also capture information about next-nearest-neighbor pairs and cross-loop (Buck, 1992; Weitzner et al., 2015) pairs, both for the stem (or "torso;" see Methods) and the entire CDR3, as well as about consecutive three-amino-acid motifs ($n=325$ non-length parameters for descriptor models vs. 420 for amino-acid models, including the 20 single-amino-acid biases). We therefore first compared 420-parameter amino acid models against 325-parameter descriptor models that fit this additional information.

We expected the descriptor models to outperform these amino-acid models, which, unlike our benchmark amino-acid models, did not fit next-nearest-neighbor pairs, reflecting the utility of additional positional constraints for defining CDR3s. We found that descriptor models outperformed amino-acid models handily, with scores of 85.6% to 14.4% for IgH (Fig. 3c, left) and 86.9% to 13.1% for TCR β (Fig. 3c, right). The median sequence in the test set was 217- and 82-fold more likely to have been produced by the descriptor model for IgH and TCR β , respectively. More remarkably, descriptor models also outperformed our benchmark amino-acid models, even though the descriptor models had less than half the parameters (820 vs. 325 non-length parameters), by almost the same margin for IgH, 80.7% to 19.3% (Fig. 3d, left), but by much less for TCR β , at 54.6% to 45.4% (Fig. 3d, right), leaving the main advantages in this case being coverage and training time. In contrast to the amino-acid models, the familiar start and end motifs (see Discussion) had already been learned in just a few iterations/minutes, requiring just a few hundred sample sequences on which to learn.

Arora et al. (2019) Repertoire-Based Diagnostics Using Statistical Biophysics

143 Test 3: Classification. Finally, we sought to test the utility of descriptor models in distinguishing
144 between states of health. As proof of principle, we fit descriptor models on 31 before-and-after
145 IgG⁺ repertoires (including three replicates) from 14 healthy human volunteers who were admin-
146 istered a seasonal trivalent influenza vaccine (Vollmers et al., 2013). We had previously shown
147 that vaccination leads to prominent changes in both repertoires' raw and functional diversity
148 (Arora et al., 2018), but withheld diversity measurements from the present study in order to test
149 the discriminatory power of the models in the absence of that additional information. Using strat-
150 ified 3-fold cross-validation, we found that descriptor models distinguished between pre- and
151 post-vaccination pairs with median AUC of 0.97 ($p=4\times 10^{-3}$; Fig. 4). It is worth noting that apply-
152 ing PCA to the models to reduce them to two dimensions, failed to distinguish between day 0
153 and day 7, consistent with a lack of necessity for directions of greatest variance to correlate with
154 differences in states of health.

155 Discussion

156 MaxEnt is a powerful method for modeling highly complex systems such as IgH and TCR β rep-
157 ertoires but exhibits practical limitations related to speed and dimensionality when fit on amino
158 acids using only the standard 20-letter alphabet. Here we demonstrate significant advantages
159 by fitting on biophysical descriptors. We show that appropriate descriptor models can capture
160 more of the information in the repertoire with fewer parameters, and that they can successfully
161 classify health-based states with high accuracy, using the IgG⁺ B-cell response to influenza
162 vaccination as proof of principle.

163 A key finding was that descriptor models outperformed amino acid models only once additional
164 positional information was included; when fit on similar positional information—single/overall
165 frequencies and nearest- and next-nearest neighbors—amino-acid models performed better.
166 This finding raises the question of what the relative contributions are of the additional types of
167 positional information fit by the winning descriptor models. There were three additional types of
168 positional information beyond nearest-neighbors: parameters for the stem, cross-pairs, and tri-
169 ples. We chose to fit the stem explicitly because the first and last few amino acids in CDR3s of
170 both IgH and TCR β are stereotypical, almost canonically beginning with a cysteine (excluded in
171 some definitions), followed by a hydroxylic or small aliphatic amino acid (most often glycine, al-
172 anine, or threonine) at the second position, and a basic amino acid (arginine/lysine) at the third
173 position and ending with a methionine or phenylalanine, followed by an aspartate, then valine or
174 tyrosine, and finally tryptophan for IgH, and starting with cysteine, alanine, and a pair of hydrox-

Arora et al. (2019) Repertoire-Based Diagnostics Using Statistical Biophysics

175 ylic or basic amino acids and often ending with glutamate, a variable amino acid, and two aro-
 176 matic amino acids for TCR β . These amino acids are important in establishing the stem-loop (or
 177 “torso-head”) configuration of CDR3s (Buck, 1992; Weitzner et al., 2015). In IgH and TCR β the
 178 stem is most often encoded by the end of the V gene segment and start of the J, not by the
 179 highly variable D gene segment and adjacent non-templated nucleotides (Lefranc et al., 1999);
 180 fitting the stem may be allowing the remaining parameters to better fit the more variable region.
 181 Fitting cross pairs—the product of descriptor values at the first and last amino acids, second
 182 and second-to-last, etc.—was similarly inspired by CDR3s’ stem-loop architecture and may be
 183 having a similar benefit. Amino-acid triples are important parts of binding motifs and have been
 184 shown to have discriminatory power in IgH in model systems (Sun et al., 2017); it is reasonable
 185 that the biophysical patterns they represent also add resolving/discriminatory power. A system-
 186 atic dissection of these contributions is left for future work.

187 A further finding was the disparity between the performance of the top descriptor model for IgH,
 188 relative to the benchmark amino-acid model, vs. that for TCR β : the descriptor model scored
 189 80.7% for IgH vs. only 54.6% for TCR β . A value over 50% indicates that the descriptor model is
 190 capturing more information than the amino-acid model, but in the case of TCR β , the benefit was
 191 modest. We considered three possible explanations. First, it is possible that both models cap-
 192 tured substantially all of the information present in the training set; however, had this been the
 193 case, the models’ final samples would likely have been nearly identical to the training set, and
 194 they were not. Second, the additional information in the TCR β repertoire may not be well cap-
 195 tured by the additional positional relationships fit by these models (stem, cross-pairs, triples),
 196 but may reside instead in some other relationship(s). Third, the modest benefit may mean that
 197 there are isolated (i.e. discontinuous) probability densities in this training set, which the Markov
 198 chain used to generate samples (Fig. 1c) has difficulty navigating (van Ravenzwaaij et al., 2018).
 199 If so, it may be that somatic hypermutation in the IgH CDR3s bridges probability densities in IgH
 200 repertoires that in TCR β repertoires, which lack somatic hypermutation, remain separate. Con-
 201 versely, the greater improvement noted for IgH may reflect descriptor models’ ability to detect
 202 biophysical similarities among these related sequences, which may be less prominent in TCR β
 203 repertoires but simultaneously difficult to capture in amino-acid models.

204 The success of descriptor models in correctly discriminating between pre- and post-influenza
 205 vaccination suggests potential medical applications. We note that vaccination, like many immu-
 206 nological perturbations, results in systems- as well as sequence-level changes; for example,
 207 changes in immunological/repertoire diversity (Jiang et al., 2013; Vollmers et al., 2013). We pre-

Arora et al. (2019) Repertoire-Based Diagnostics Using Statistical Biophysics

viously showed that the combination of raw and functional diversity, measured with various frequency weightings, can discriminate between pre- and post-vaccination sample pairs with high accuracy, likely in part by detecting clonal expansion with selection (Arora et al., 2018). However, changes in diversity, while potentially useful as part of a screening test, are not sufficiently specific to serve as a general diagnostic modality. The present study shows that even without the powerful discriminatory information that diversity adds, descriptor models are capable of highly sensitive and specific diagnostic discrimination, with high AUC and low p-value from small numbers of subjects and samples. The relatively small number of parameters and these parameters' relatively straightforward interpretability (compared to, for example, parameters in deep-learning models) suggest that leveraging the statistical biophysics of repertoires' amino acid composition is a promising direction for dissecting immune responses for diagnostic and therapeutic purposes. This method is extensible to more or all of IgH or TCR β , to the complementary chain (IgL/TCR α), and indeed to other proteins or biopolymers, leveraging the power of functional relationships to shrink alphabets while increasing their information density.

Acknowledgements

This work used the resources allocated via Jetstream cloud service (allocation ID: TG-BIO170094) of Extreme Science and Engineering Discovery Environment (XSEDE), which is supported by National Science Foundation grant number ACI-1548562. Research Computing Group at the High Performance Computing Cluster at Harvard Medical School. The authors would like to thank Dr. Mohammed Al-Quraishi for helpful discussion.

Methods

Descriptors. Twenty-six biophysical measurements were previously made on a set of 87 amino acids, which included the standard 20 (Sandberg et al., 1998). We filtered out non-standard amino acids and applied PCA to the standard 20 amino acids (using Python; sklearn.decomposition.PCA library). The top five PCs, which together explained 92% of the observed variance, were each normalized to a mean of 0 and maximum range [-1, 1] and used as biophysical descriptors.

Data. IgG (Vollmers et al., 2013), memory IgH (DeWitt et al., 2016), and TCR β (Emerson et al., 2017) CDR3 repertoires were obtained and processed as previously described (Arora et al., 2018). For each dataset in Tests 1 and 2, 500,000 sequences were chosen at random and split 90:10 into training and test sets; for Test 3 all sequences were used.

Arora et al. (2019) Repertoire-Based Diagnostics Using Statistical Biophysics

239 Models. MaxEnt models were trained on features' expectation values, with one parameter per
 240 feature (the bias). For Tests 1 and 2, amino-acid models were trained on the observed frequen-
 241 cies of the single amino acids ($n=20$ parameters) and nearest- and next-nearest-neighbor ami-
 242 no-acid pairs ($n=20^2 \times 2 = 800$) and the frequencies of CDR3 lengths ($n=38$ for IgH and 26 for
 243 TCR β), following previous reports (Mora et al., 2010). Descriptor models were trained on the
 244 frequencies of the single amino acids ($n=20$ parameters), the product of each pair of descriptors
 245 at different positions ($n=5^2=25$ per set), lengths, and, as indicated, on amino-acid frequencies
 246 for the first- and last-four amino acids (roughly corresponding to the CDR3 stem or "torso"
 247 (North et al., 2011; Finn et al., 2016); $n=20$), the product of each (non-redundant) pair of de-
 248 scriptors at the same position ($n=(5 \times 4)/2 = 10$), the product of each pair of descriptors for the
 249 stem ($n=25$ per set), and the product of descriptors at each three ($n=5^3=125$). For Test 3, mod-
 250 els were trained on the expectation values of each descriptor for the stem ($n=5$) and full-length
 251 CDR3 ($n=5$), pairs of descriptors at the same position, cross-loop pairs, and nearest- and next-
 252 nearest-neighbor pairs for the full-length CDR3 and the stem ($n=25$ per set), anchoring se-
 253 quences with an initial cysteine and terminal tryptophan for speed.

254 Fitting was performed using Metropolis-Hastings Markov-chain Monte Carlo sampling with the
 255 acceptance criterion

$$A(\sigma', \sigma) = \min \left(\frac{p(\sigma')}{p(\sigma)} \frac{g(\sigma'|\sigma)}{g(\sigma|\sigma')}, 1 \right)$$

256 where σ is the original sequence and σ' is proposed according to the proposal distribution
 257 $g(\sigma|\sigma')$, updating biases via gradient descent using an adaptive step size, using an adaptive
 258 burn-in period and autocorrelation time, and a time limit of 24 hours/fit as a stopping condition.
 259 Each model was trained for 24 hours on 44 parallel CPUs using the National Science Founda-
 260 tion's high-performance supercomputing cluster, XSEDE (Towns et al., 2014). To avoid overfit-
 261 ting, we prohibited sample size from exceeding the size of the training set.

262 Probabilities. The probability of a sequence σ according to a MaxEnt model M was calculated as

$$p(\sigma|M) = \frac{1}{Z} e^{-E_{\sigma|M}}$$

Arora et al. (2019) Repertoire-Based Diagnostics Using Statistical Biophysics

where $E_{\sigma|M}$ is the energy of σ and the normalization constant $Z = \sum_{\sigma} e^{-E_{\sigma|M}}$ was estimated via bridge sampling (Meng and Wong, 1996; Gelman and Meng, 1998) using Harvard Medical School's high-performance computing cluster.

Classification (Test 3). We fit descriptor models on each of the day 0/day 7 before-and-after IgG⁺ repertoire pairs ($n=31$: $n=17$ from day 0, including replicates, and $n=14$ day 7) from the influenza vaccination dataset (Vollmers et al., 2013; Arora et al., 2018) and used a support-vector-machine (SVM) on the final models for classification (excluding length biases, which interact with the normalization constant), using the median area under the receiver-operator-characteristic curve (AUC/ROC) as the quality measure (taken over $n=10,000$ repeats; mean preferred over median given the observed highly skew AUC distributions expected from strong performance with outliers; Fig. 4 top inset), with stratified k -fold cross-validation (without over-sampling; 17 vs. 14 was considered sufficiently balanced, but see null-model comparison below) to avoid overfitting (for $k=2, 3, 5$, and 10 to confirm robustness) and comparison to the AUC of randomly relabeled data as a null model (also $n=10,000$ repeats) to assess statistical significance. Mann-Whitney U p -value was calculated to test that the two AUC distributions were different. The significance of the AUC was understood as the probability that it could arise from a random classifier by chance; the p -value for significance of the AUC was therefore calculated as the fraction of the area under the null-model distribution to the right of the AUC. Histograms were plotted. All analyses were performed using Python's numpy and scipy libraries.

Figure Legends

Figure 1. MaxEnt Based on Amino Acids' Biophysical Properties. (a) Amino acids as vectors, shown here as a heatmap, in a 5-dimensional descriptor space. (b) Amino acids with similar properties lie near to each other in descriptor space. These similarities can be visualized by calculating all pairwise Euclidean distances of the amino acids in descriptor space, constructing a (complete, K_{20}) network with the amino acids as nodes and the distances as weighted edges, and then for clarity keeping only edges with weights ≤ 1.1 . For example, aspartate (D) and glutamate (E) (red boxes in (a)) lie near to each other in descriptor space, illustrated by their similar pattern in the heatmap (with prominent differences only in the dimension corresponding to descriptor 4), and so are adjacent in the network. Amino acids are colored according to a familiar groupings (basic, aliphatic, etc.) to demonstrate that their configuration in descriptor space agrees with these groupings. (c) Data preparation and model training. Repertoires were first split into training and test sets, and the features of the training set measured. Models were

Arora et al. (2019) Repertoire-Based Diagnostics Using Statistical Biophysics

trained through iterative sampling, comparing of sample and observed features, and updating biases. **(d)** Example using a highly simplified toy repertoire consisting of a training set of two unique 3-amino-acid sequences, CTR and DVW (common in stems of IgH CDR3s). The models learn constraints that distinguish the training set from random 3-mers. For amino-acid models, constraints are the frequencies of letters; for descriptor models, constraints are expectation values of descriptors and descriptor products at given positions (here, nearest-neighbor pairs). Model output is shown in the last row. The descriptor model has learned the pattern of biophysical relationships, such that sequences that are biophysically similar to sequences in the training set also appear in the sample, albeit at lower frequency than the sequences in the training set.

Figure 2. Training and normalization. Distance (root-mean-squared error, RMSE) between training data and model sample as a function of iterations of model training. Data shown is for all models in Tests 1 and 2. Insets, bridge sampling for representative fits showing overlap between model- (blue) and randomly sampled sequences (gray).

Figure 3. Comparison of amino-acid vs. descriptor models. Head-to-head tests on IgH (left) and TCR β (right) repertoires; the better performer is shaded green. **(a)** Validation comparison of models of IgH vs. TCR β repertoires; IgH models strongly prefer IgH sequences (yellow) and TCR β models strongly prefer TCR β sequences (red; results shown are for the 325-parameter descriptor models). **(b)-(d)** Comparisons of an amino-acid model to a descriptor model, both trained/tested on the same training/test set. Density to the left of the vertical dashed line represents sequences for which the amino-acid model gave the higher probability; density to the right (filled) represents higher probability per the descriptor model. Vertical red lines denote medians of the probability densities. **(b)** Test 1: models fitting similar positional information (single positions plus nearest- and next-nearest neighbors); amino-acid models perform better. **(c)** Test 2: models fitting similar numbers of parameters (420 non-length parameters for the amino-acid model vs. 325 for the descriptor model); descriptor models perform better. **(d)** Test 2, continued: amino-acid benchmark model (820 parameters; nearest- and next-nearest neighbors) vs. the descriptor model in (c); descriptor models perform better.

Figure 4. Classification of pre- vs. post-flu vaccination in human subjects. Shown is the median AUC (red) for 10,000 training-test splits using stratified 3-fold cross-validation of an SVM on 31 pre- and post-vaccination samples from the same subjects. Insets show the distributions of AUCs from all 10,000 splits of the real data (blue) and from 10,000 splits in which the data was randomly relabeled, to measure the probability that the median performance could have been

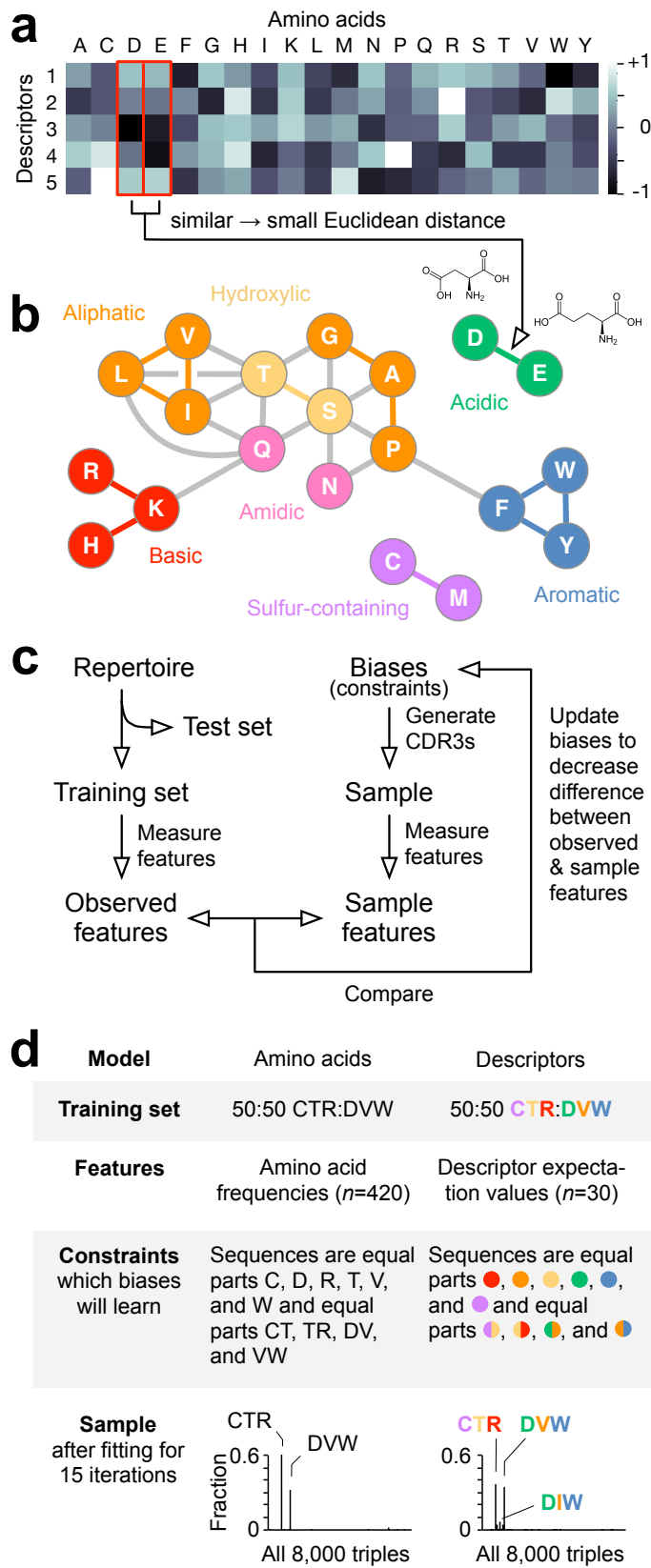
Arora et al. (2019) Repertoire-Based Diagnostics Using Statistical Biophysics

327 the result of chance (gray). Red, median. The p-value is the area in the random-relabeling dis-
328 tribution to the right of the median.

Arora et al. (2019) Repertoire-Based Diagnostics Using Statistical Biophysics

329

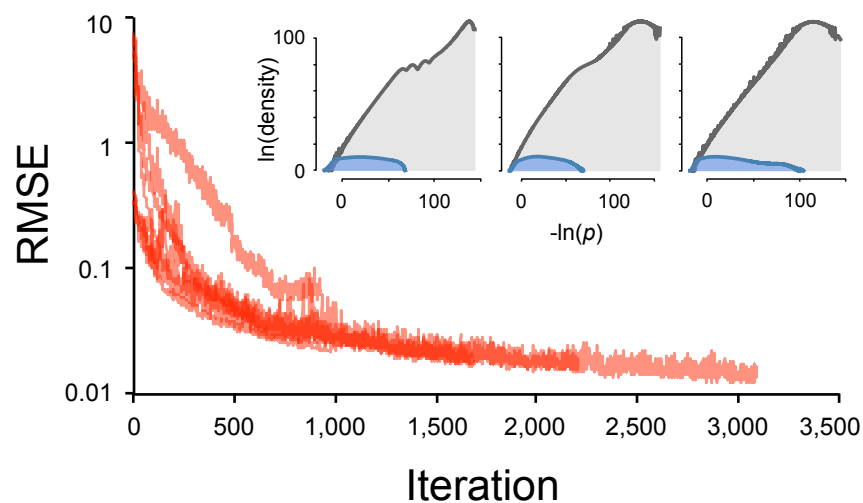
Figure 1



330

Arora et al. (2019) Repertoire-Based Diagnostics Using Statistical Biophysics

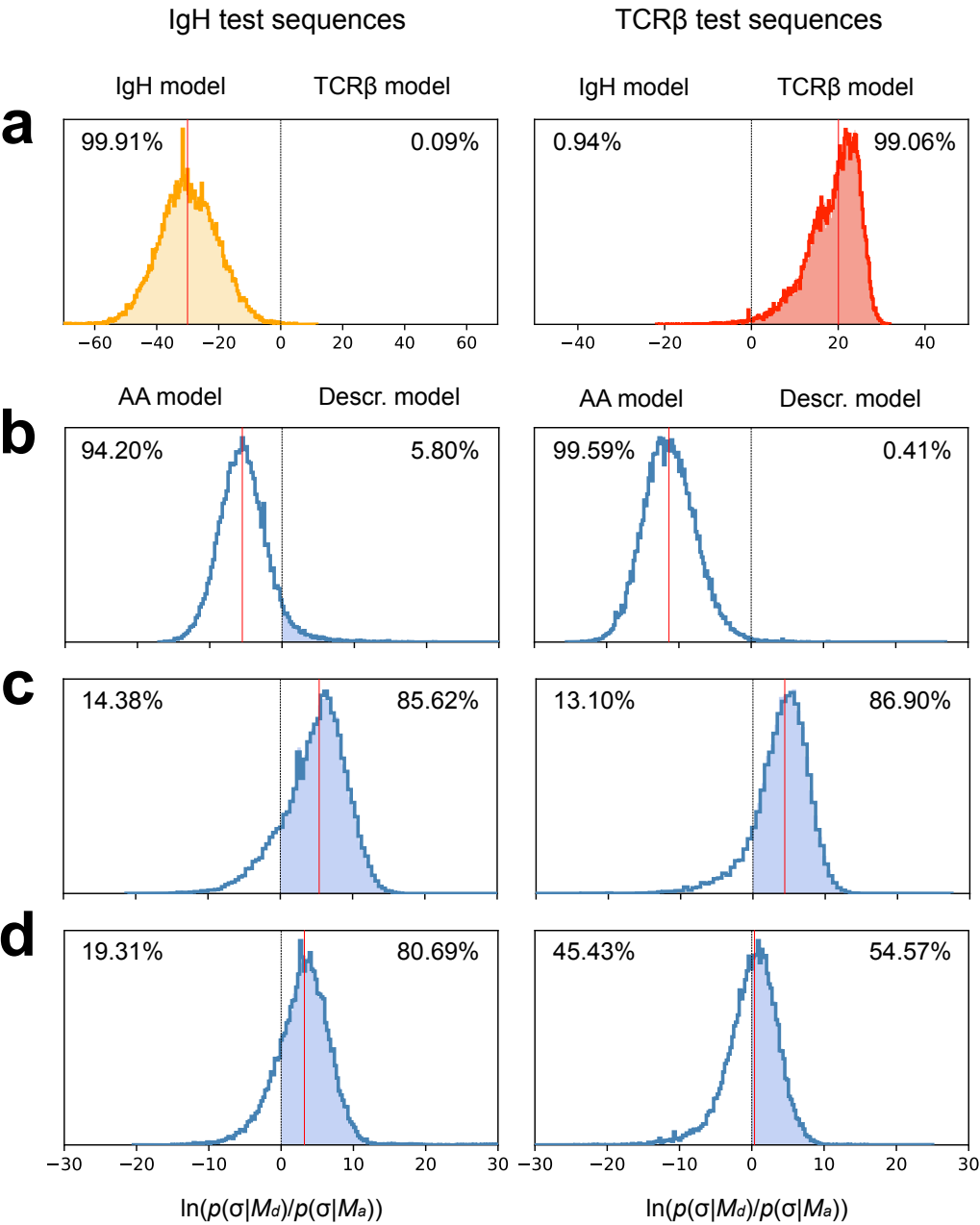
Figure 2



Arora et al. (2019) Repertoire-Based Diagnostics Using Statistical Biophysics

333

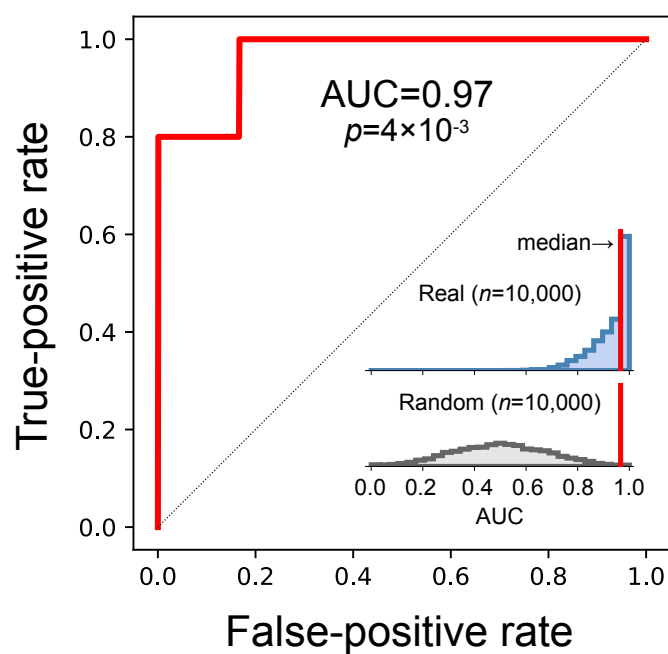
Figure 3



334

Arora et al. (2019) Repertoire-Based Diagnostics Using Statistical Biophysics

Figure 4



Arora et al. (2019) Repertoire-Based Diagnostics Using Statistical Biophysics

337 References

- 338 Arnaout, R., Lee, W., Cahill, P., Honan, T., Sparrow, T., Weiland, M., Nusbaum, C., Rajewsky,
339 K., and Koralov, S.B. (2011). High-Resolution Description of Antibody Heavy-Chain Repertoires
340 in Humans. *PLOS ONE* 6, e22365.
- 341 Arora, R., Burke, H.M., and Arnaout, R. (2018). Immunological Diversity with Similarity. *BioRxiv*
342 483131.
- 343 Bialek, W., and Ranganathan, R. (2007). Rediscovering the power of pairwise interactions.
344 *ArXiv:0712.4397 [q-Bio]*.
- 345 Buck, C.A. (1992). Immunoglobulin superfamily: structure, function and relationship to other re-
346 ceptor molecules. *Semin. Cell Biol.* 3, 179–188.
- 347 DeWitt, W.S., Lindau, P., Snyder, T.M., Sherwood, A.M., Vignali, M., Carlson, C.S., Greenberg,
348 P.D., Duerkopp, N., Emerson, R.O., and Robins, H.S. (2016). A Public Database of Memory
349 and Naive B-Cell Receptor Sequences. *PLoS ONE* 11, e0160853.
- 350 Dosztányi, Z., and Torda, A.E. (2001). Amino acid similarity matrices based on force fields. *Bio-*
351 *informatics* 17, 686–699.
- 352 Emerson, R.O., DeWitt, W.S., Vignali, M., Gravley, J., Hu, J.K., Osborne, E.J., Desmarais, C.,
353 Klinger, M., Carlson, C.S., Hansen, J.A., et al. (2017). Immunosequencing identifies signatures
354 of cytomegalovirus exposure history and HLA-mediated effects on the T cell repertoire. *Nat.*
355 *Genet.* 49, 659–665.
- 356 Finn, J.A., Koehler Leman, J., Willis, J.R., Cisneros, A., Crowe, J.E., and Meiler, J. (2016). Im-
357 proving Loop Modeling of the Antibody Complementarity-Determining Region 3 Using
358 Knowledge-Based Restraints. *PLoS ONE* 11, e0154811.
- 359 Gelman, A., and Meng, X.-L. (1998). Simulating Normalizing Constants: From Importance Sam-
360 pling to Bridge Sampling to Path Sampling. *Statistical Science* 13, 163–185.
- 361 Greiff, V., Weber, C.R., Palme, J., Bodenhofer, U., Miho, E., Menzel, U., and Reddy, S.T. (2017).
362 Learning the High-Dimensional Immunogenomic Features That Predict Public and Private Anti-
363 body Repertoires. *J. Immunol.* 199, 2985–2997.

Arora et al. (2019) Repertoire-Based Diagnostics Using Statistical Biophysics

- 364 Hellberg, S., Sjöström, M., Skagerberg, B., and Wold, S. (1987). Peptide quantitative structure-
365 activity relationships, a multivariate approach. *J. Med. Chem.* *30*, 1126–1135.
- 366 Jiang, N., He, J., Weinstein, J.A., Penland, L., Sasaki, S., He, X.-S., Dekker, C.L., Zheng, N.-Y.,
367 Huang, M., Sullivan, M., et al. (2013). Lineage structure of the human antibody repertoire in re-
368 sponse to influenza vaccination. *Sci Transl Med* *5*, 171ra19.
- 369 Kaplinsky, J., Li, A., Sun, A., Coffre, M., Koralov, S.B., and Arnaout, R. (2014). Antibody reper-
370 toire deep sequencing reveals antigen-independent selection in maturing B cells. *Proc. Natl.*
371 *Acad. Sci. U.S.A.* *111*, E2622-2629.
- 372 Kim, S., Thiessen, P.A., Bolton, E.E., Chen, J., Fu, G., Gindulyte, A., Han, L., He, J., He, S.,
373 Shoemaker, B.A., et al. (2016). PubChem Substance and Compound databases. *Nucleic Acids*
374 *Res.* *44*, D1202-1213.
- 375 Lefranc, M.P., Giudicelli, V., Ginestoux, C., Bodmer, J., Müller, W., Bontrop, R., Lemaitre, M.,
376 Malik, A., Barbié, V., and Chaume, D. (1999). IMGT, the international ImMunoGeneTics data-
377 base. *Nucleic Acids Res.* *27*, 209–212.
- 378 Marks, D.S., Colwell, L.J., Sheridan, R., Hopf, T.A., Pagnani, A., Zecchina, R., and Sander, C.
379 (2011). Protein 3D Structure Computed from Evolutionary Sequence Variation. *PLOS ONE* *6*,
380 e28766.
- 381 Meng, X.-L., and Wong, W.H. (1996). Simulating ratios of normalizing constants via a simple
382 identity: A theoretical exploration. *Statistica Sinica* *6*, 831–860.
- 383 Mora, T., Walczak, A.M., Bialek, W., and Callan, C.G. (2010). Maximum entropy models for an-
384 tibody diversity. *PNAS* *107*, 5405–5410.
- 385 North, B., Lehmann, A., and Dunbrack, R.L. (2011). A new clustering of antibody CDR loop con-
386 formations. *J. Mol. Biol.* *406*, 228–256.
- 387 Ostmeier, J., Christley, S., Rounds, W.H., Toby, I., Greenberg, B.M., Monson, N.L., and Cowell,
388 L.G. (2017). Statistical classifiers for diagnosing disease from immune repertoires: a case study
389 using multiple sclerosis. *BMC Bioinformatics* *18*, 401.

Arora et al. (2019) Repertoire-Based Diagnostics Using Statistical Biophysics

- 390 Ostmeyer, J., Christley, S., Toby, I.T., and Cowell, L.G. (2019). Biophysicochemical motifs in T
391 cell receptor sequences distinguish repertoires from tumor-infiltrating lymphocytes and adjacent
392 healthy tissue. *Cancer Res canres*.2292.2018.
- 393 Parameswaran, P., Liu, Y., Roskin, K.M., Jackson, K.K.L., Dixit, V.P., Lee, J.-Y., Artiles, K.L.,
394 Zompi, S., Vargas, M.J., Simen, B.B., et al. (2013). Convergent antibody signatures in human
395 dengue. *Cell Host Microbe* 13, 691–700.
- 396 van Ravenzwaaij, D., Cassey, P., and Brown, S.D. (2018). A simple introduction to Markov
397 Chain Monte–Carlo sampling. *Psychon Bull Rev* 25, 143–154.
- 398 Robins, H.S., Srivastava, S.K., Campregher, P.V., Turtle, C.J., Andriesen, J., Riddell, S.R.,
399 Carlson, C.S., and Warren, E.H. (2010). Overlap and Effective Size of the Human CD8+ T Cell
400 Receptor Repertoire. *Science Translational Medicine* 2, 47ra64-47ra64.
- 401 Russ, W.P., Lowery, D.M., Mishra, P., Yaffe, M.B., and Ranganathan, R. (2005). Natural-like
402 function in artificial WW domains. *Nature* 437, 579–583.
- 403 Sandberg, M., Eriksson, L., Jonsson, J., Sjöström, M., and Wold, S. (1998). New chemical de-
404 scriptors relevant for the design of biologically active peptides. A multivariate characterization of
405 87 amino acids. *J. Med. Chem.* 41, 2481–2491.
- 406 Seno, F., Trovato, A., Banavar, J.R., and Maritan, A. (2008). Maximum entropy approach for
407 deducing amino Acid interactions in proteins. *Phys. Rev. Lett.* 100, 078102.
- 408 Stephenson, J.D., and Freeland, S.J. (2013). Unearthing the root of amino acid similarity. *J. Mol.*
409 *Evol.* 77, 159–169.
- 410 Sun, Y., Best, K., Cinelli, M., Heather, J.M., Reich-Zeliger, S., Shifrut, E., Friedman, N., Shawe-
411 Taylor, J., and Chain, B. (2017). Specificity, Privacy, and Degeneracy in the CD4 T Cell Recep-
412 tor Repertoire Following Immunization. *Front. Immunol.* 8.
- 413 Susko, E., and Roger, A.J. (2007). On reduced amino acid alphabets for phylogenetic inference.
414 *Mol. Biol. Evol.* 24, 2139–2150.
- 415 Towns, J., Cockerill, T., Dahan, M., Foster, I., Gaither, K., Grimshaw, A., Hazlewood, V., Lath-
416 rop, S., Lifka, D., Peterson, G.D., et al. (2014). XSEDE: Accelerating Scientific Discovery. *Com-*
417 *puting in Science & Engineering* 16, 62–74.

Arora et al. (2019) Repertoire-Based Diagnostics Using Statistical Biophysics

- 418 Vollmers, C., Sit, R.V., Weinstein, J.A., Dekker, C.L., and Quake, S.R. (2013). Genetic meas-
419 urement of memory B-cell recall using antibody repertoire sequencing. *Proc. Natl. Acad. Sci.*
420 *U.S.A.* *110*, 13463–13468.
- 421 Walter, K.U., Vamvaca, K., and Hilvert, D. (2005). An active enzyme constructed from a 9-
422 amino acid alphabet. *J. Biol. Chem.* *280*, 37742–37746.
- 423 Weitzner, B.D., Dunbrack, R.L., and Gray, J.J. (2015). The origin of CDR H3 structural diversity.
424 *Structure* *23*, 302–311.
- 425 Yeo, G., and Burge, C.B. (2004). Maximum entropy modeling of short sequence motifs with ap-
426 plications to RNA splicing signals. *J. Comput. Biol.* *11*, 377–394.