

An antibody developability triaging pipeline exploiting protein language models

James Sweet-Jones & Andrew C.R. Martin

To cite this article: James Sweet-Jones & Andrew C.R. Martin (2025) An antibody developability triaging pipeline exploiting protein language models, mAbs, 17:1, 2472009, DOI: 10.1080/19420862.2025.2472009

To link to this article: <https://doi.org/10.1080/19420862.2025.2472009>



© 2025 The Author(s). Published with license by Taylor & Francis Group, LLC.



View supplementary material [↗](#)



Published online: 04 Mar 2025.



Submit your article to this journal [↗](#)



View related articles [↗](#)



View Crossmark data [↗](#)

An antibody developability triaging pipeline exploiting protein language models

James Sweet-Jones and Andrew C.R. Martin 

Institute of Structural and Molecular Biology, Division of Biosciences, University College London, London, UK

ABSTRACT

Therapeutic monoclonal antibodies (mAbs) are a successful class of biologic drugs that are frequently selected from phage display libraries and transgenic mice that produce fully human antibodies. However, binding affinity to the correct epitope is necessary, but not sufficient, for a mAb to have therapeutic potential. Sequence and structural features affect the developability of an antibody, which influences its ability to be produced at scale and enter trials, or can cause late-stage failures. Using data on paired human antibody sequences, we introduce a pipeline using a machine learning approach that exploits protein language models to identify antibodies which cluster with antibodies that have entered the clinic and are therefore expected to have developability features similar to clinically acceptable antibodies, and triage out those without these features. We propose this pipeline as a useful tool in candidate selection from large libraries, reducing the cost of exploration of the antibody space, and pursuing new therapeutics.

ARTICLE HISTORY

Received 31 October 2024
Revised 17 February 2025
Accepted 20 February 2025

KEYWORDS

Antibodies; developability; machine learning; prediction; protein language models

Introduction

Monoclonal antibodies (mAbs) have been shown to be a successful class of biologic drugs which have potential to treat a wide variety of diseases owing to their ability to target a specific antigen, and therefore potentially any step in a disease pathway.^{1,2} As of early 2025, at least 130 mAbs have received regulatory approval from the U.S. Food and Drug Administration or the European Medicines Agency (db.antibodysociety.org/) with at least 42 being considered as ‘fully-human’, either from transgenic mice, phage display libraries, or cloned from recovering patients.^{3–5} The annual growth of this sector has increased by between 20% and 30% per year^{6,7} and is likely to continue to grow as interest increases in the use of antibodies to target previously undruggable targets.⁸ Despite this, throughout the clinical pipeline for the development of new mAbs, there is a high risk of failure, causing costly discontinuation from trials.⁹

Simultaneously, efforts in single-cell sequencing techniques have been applied to understand how the antibody repertoire functions and changes over time at the level of single B cells.^{10–13} This has given researchers the ability to generate dense digital libraries of paired variable heavy (V_H) and variable light (V_L) human antibody sequences that vastly outnumber previous databases resulting from sequence or structural data (KabatMan,¹⁴ IMGT,¹⁵ SABDAb¹⁶ AbDb¹⁷ and EMBLIG (aby.bank.org/emblig/)). Online repositories including the Observed Antibody Space (OAS),¹⁸ cAb-Rep¹⁹ and BRepertoire²⁰ allow researchers access to these resources.

With the generation of these *in silico* databases, efforts to develop screening statistics to identify sequences with physical characteristics similar to approved therapeutics has become a driver in the field. Usually, these have been based on antibody

developability, which is loosely defined as an antibody’s intrinsic ability to be produced on an industrial scale, to maintain reasonable stability in long-term storage and in patients, and to be safely tolerated by the patient.^{21,22} Such considerations have now become important in the early stages of drug screening to select the best quality candidates and avoid costly late-stage failures.¹ Furthermore, developability is important, but does not guarantee success in clinical trials, where candidates may face discontinuation for safety or efficacy reasons. Identifying factors important in determining success in clinical trials has also eluded researchers.

Physicochemical features, including surface charged patches, surface hydrophobic patches, low thermostability, and post-translational modification sites that introduce heterogeneity, have become associated with poor antibody developability.²³ Those features that compromise the stability of the antibody can cause unfolding, increase the propensity to aggregate in solution and can increase immunogenicity.^{24,25} At the lead candidate stage, well-defined experimental assays for measurement are important in the selection of a final lead.^{26,27} However, it has become useful to predict these features at an earlier stage using computational means. To this end, sequence-based statistics have been developed based on these features and are available for use in drug discovery pipelines, including the Developability Index,^{28,29} AbPred,³⁰ and, more recently, the Therapeutic Antibody Profiler (TAP)³¹ and Therapeutic Antibody Developability Analysis (TA-DA Score).³² However, these tools can fall short in identifying leads from large libraries of data, requiring computationally expensive 3D modeling, or only taking one antibody at a time, which is usually expected already to be a potential lead candidate.

CONTACT Andrew C.R. Martin  andrew@bioinf.org.uk; andrew.martin@ucl.ac.uk  Institute of Structural and Molecular Biology, Division of Biosciences, University College London, London, UK

 Supplemental data for this article can be accessed online at <https://doi.org/10.1080/19420862.2025.2472009>

© 2025 The Author(s). Published with license by Taylor & Francis Group, LLC.

This is an Open Access article distributed under the terms of the Creative Commons Attribution License (<http://creativecommons.org/licenses/by/4.0/>), which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited. The terms on which this article has been published allow the posting of the Accepted Manuscript in a repository by the author(s) or with their consent.

In order to take advantage of the wealth of data now available, the field has also turned to machine learning as a new avenue of exploration.^{1,33,34} For protein sequences to be suitable inputs for machine learning problems, it is necessary to encode them numerically. Previously, this has been done by using evolutionary or physicochemical and structural features,^{35–37} and simple regression models to identify features of high importance, or to predict features from the sequence as done in AbPred.³⁰ Negron et al.³² expanded on this work and identified previously mentioned characteristics, including hydrophobicity (assessed by hydrophobic interaction chromatography), thermostability (T_m , assessed by differential scanning fluorimetry) and aggregation (assessed by cross-interaction chromatography) that were associated with the identification of clinically acceptable mAbs. Furthermore, this work has demonstrated an ability to separate clinical antibody sequences from antibody repertoires and to assign a developability score based on these features as part of their TA-DA score.

Many studies, including those described above^{28–32} and others³⁸ as well as reviews,^{39–41} have described the importance of predicting developability and most of these approaches rely on the assumption that clinical antibodies (i.e., approved, discontinued and in-development mAbs) have a range of properties related to developability such that novel antibodies with similar (predicted) properties will also be clinically acceptable. This is not to say that antibodies with very different properties will necessarily fail in the clinic, but that such approaches allow one to focus on antibodies most likely to succeed. The need to exploit ‘big data’ and artificial intelligence in the development of biologics has been discussed by Fernández-Quintero et al.³⁹ and Narayanan et al.⁴² A newer method of encoding protein sequences is to use ‘protein language models’.⁴³ These are deep learning encoders trained on the relationships between residues in a sequence using millions of sequences. The results give dense numerical representations of sequences that may then be used as training data for machine learning models and, over the last few years, have revolutionized predictive methods in all areas of bioinformatics (see Lin et al.,⁴⁴ for example). Their power comes from their ability to encode more information, including less obvious features and combinatorial or multi-factor features (e.g., from interaction of amino acids).

In this study, we hypothesized that, rather than directly predicting physical properties related to developability, antibodies with developable traits may be selected by encoding them using protein language models and comparing the encoded antibodies with encoded sequences of current clinical mAbs. Thus, our approach is, in principle, similar to many

other approaches (including TAP), that look for similarity in properties to clinical antibodies, but, unlike these methods, we do not explicitly predict developability features, but instead exploit the power of protein language models for encoding antibody sequences.

Our goal is then to build a high-throughput triaging pipeline exploiting preliminary simple physicochemical screening followed by machine learning using protein language models which may be used to select antibodies most likely to have good developability characteristics from large libraries.

Results

Simple physicochemical properties of clinical and library antibodies

As a first step, we looked at using physicochemical properties to attempt to identify antibodies with clinically acceptable properties in a set of library antibodies. The aim was to see whether the clinical mAbs have a restricted distribution of these properties compared with antibodies from a library, similar to the approach used by Raybould *et al.*,³¹ except here, we use only sequence statistics that can be calculated quickly without high computational expense.

A dataset was collated consisting of paired V_H and V_L sequences of clinical stage human mAbs ($n=144$) from the October 2021 release of TheraSabDab³ marked as ‘Whole mAb’ (Supplementary Table S1) and 10,000 paired sequences randomly selected from the OAS online repertoire repository (accessed January 2022)¹⁸ (Supplementary Table S2). We refer to this set of sequences from OAS as our ‘library’.

Physicochemical properties, including predicted ΔG of unfolding,⁴⁵ iso-electric point (pI)⁴⁶ and CDR-H3 loop length⁴⁷ were calculated. Using a Mann-Whitney U test, we observed that there were statistical differences in the CDR-H3 length and in the predicted ΔG of unfolding for concatenated V_H and V_L chains between therapeutic and library antibodies (Table 1 and Supplementary Table S3). While this demonstrates a difference between human repertoire antibodies and what is found in the clinical mAb dataset, the mean values are relatively similar in the two datasets, making it difficult to use this as an approach to identify antibodies with clinically acceptable developability characteristics, although it can be used to reject clear outliers.

Table 1. Means and standard deviation of sequence-calculated physicochemical properties for fully human mAb therapeutics ($n=144$) and repertoire human antibodies from OAS ($n=10,000$, the ‘Human Library Antibodies’). p-values were calculated using a Mann-Whitney U test.

| Feature | Human Therapeutic mAbs | Human LibraryAntibodies | p-value |
|---|------------------------|-------------------------|---------|
| CDR-H3 Loop Length | 12.1 ± 6.65 | 15.0 ± 10.54 | 0.00049 |
| $\Delta G V_H$ (kJ mol ⁻¹) | 7614 ± 3260 | 6583 ± 3441 | 0.00014 |
| $\Delta G V_L$ (kJ mol ⁻¹) | 1086 ± 2381 | 796 ± 2614 | 0.14 |
| Concatenated V_H/V_L ΔG (kJ mol ⁻¹) | 9248 ± 3896 | 7944 ± 4238 | 0.00015 |
| Mean pI of V_H/V_L | 7.9 ± 1.30 | 7.8 ± 1.24 | 0.025 |

Identifying clinical-like antibodies from repertoires using unsupervised learning

An unsupervised learning model was proposed as an approach to identify library antibodies with clinical-antibody-like properties. Just as with approaches such as TAP, we hypothesize that clinical mAbs (which have probably undergone some developability assessment) should cluster in some N -dimensional space and that repertoire antibodies with similar properties would be positioned close to the clinical mAbs. To train an unsupervised learning method, the library and clinical V_H and V_L sequences were padded according to the Chothia numbering scheme, then independently encoded with various language models: ESM,⁴⁴ AbLang,⁴⁸ Sapiens⁴⁹ and AntiBERTy.⁵⁰ The encodings generated 130,048 features per paired V_H/V_L sequence. All language models had a similar performance for this task, with AntiBERTy somewhat outperforming the other methods (data not shown).

Various unsupervised machine learning models were tested: linear Principal Component Analysis (PCA),⁵¹ kernel PCA,⁵¹ 2-dimensional (2D) 't-distributed Stochastic Neighbor Embedding' (t-SNE)⁵² and 'Uniform Manifold Approximation and Projection' (UMAP)⁵³ (Figure 1a). These algorithms demonstrate how library antibodies are positioned against clinical mAbs also encoded with the AntiBERTy language model. For the linear PCA, t-SNE or UMAP, data were arranged into discrete groups of antibodies which are dictated by V_H and V_L gene germline pairing (Figure 1b). However, kernel PCA with a radial basis kernel function ($\gamma = 500$, Supplementary Figure S1), when viewing the first two principal components, gave a useful pattern of clustering where library antibodies form a radial pattern with clinical mAbs positioned around the origin (Figure 1a). This was also true of a held-back dataset of human-derived clinical mAbs ($n=203$) named with the 2016 and 2022 naming conventions⁵⁴ in which the source infix ('-u-' for human or '-zu-' for humanized) was removed, and therefore human mAbs could not be identified using the '-umab' but not '-zumab' approach used to identify human mAbs with the earlier naming schemes (see Methods and Supplementary Table S4). These held-back antibodies were positioned close to the original dataset of human-clinical mAbs (Figure 2). This led us to conclude that repertoire antibodies which are positioned close to clinical mAbs may be likely to share the developability properties necessary and should be taken forward for potential development.

Cutoffs were then established to select the repertoire antibodies which cluster with the clinical mAbs in order to extract them. An ellipse function was used in which the principal component with the greater range for clinical mAbs was taken to be the major axis, and the lesser range as the minor axis. Z-score thresholds (the number of standard deviations away from the mean) along the two principal components of the clinical mAbs were used to select where the extremes of the ellipse should be placed. The Z-score thresholds were optimized by measuring the proportion of the clinical mAbs captured by the ellipse against the proportion of the library antibodies also captured in the same ellipse. It was expected that, since the spread of antibodies was even across the first

two principal components of the kernel PCA, roughly equal proportions of both groups would be captured. This was done using all human clinical mAbs (Figure 3a) and with only approved human mAbs (Figure 3b).

Comparing Figures 3a and 3b it can be seen that the bars for the OAS (library) antibodies are consistently lower when the Z-scores are based on the approved antibodies than they are when based on the clinical (i.e., approved, discontinued and in-development) antibodies. This indicates that the approved antibodies occupy a tighter distribution than the clinical antibodies. While it is obvious that the approved antibodies will be a subset of the clinical antibodies, it is less obvious that they will form a tighter cluster in this projection of the AntiBERTy-encoded parameter space. This led us to conclude that there may be characteristics of the approved antibodies identified by the protein language model that would allow them to be separated from the antibodies that were discontinued.

Using supervised machine learning to distinguish approved and discontinued clinical antibodies

It is evident that having suitable developability profiles alone is not sufficient for an antibody to succeed in clinical trials and the clinical dataset used to identify library antibodies with properties similar to clinical mAbs contained discontinued antibodies. There are many reasons why an antibody could fail in clinical trials. Some of these are intrinsic to the sequence (e.g., immunogenicity, developability), while others are target-specific (e.g., binding affinity, nature of the epitope, on- or off-target side-effects).^{1,9,55} It is also likely that the effectiveness threshold for a drug to be taken forward from Phase 3 trials will be higher if there are already effective drugs on the market. However, given the differences observed above, we considered it worthwhile to attempt to train a predictor that might be able to identify drug-like antibodies that are more likely to succeed in the clinic. We assembled a dataset of the V_H and V_L amino acid sequences for 115 approved and 150 discontinued antibodies from the TheraSabDab³ (Supplementary Table S5).

Unlike the comparison of human clinical mAbs and library antibodies, which had statistical differences in their physicochemical properties (Table 1), there were no statistically significant differences in any of the basic physicochemical properties between approved and discontinued antibodies using a Mann-Whitney U test. This included the G score⁵⁶ used as a method for predicting immunogenicity (Table 2). The largest quantitative difference was that the discontinued antibodies had a lower mean length for the CDR-H3 loop (Supplementary Table S6). The lack of statistically significant differences is perhaps not surprising given that both approved and discontinued antibodies will almost certainly have undergone a developability assessment and possibly optimization before entering clinical trials. It was also seen that the approved and discontinued groups had similar proportions of V_H and V_L V-gene germline pairings (Supplementary Figure S2).

As before, the V_H and V_L sequences for each antibody were padded and aligned using the Chothia numbering scheme and the sequences were then encoded with a selection of general protein (ESM⁴⁴) and antibody-specific (AntiBERTy,⁵⁰

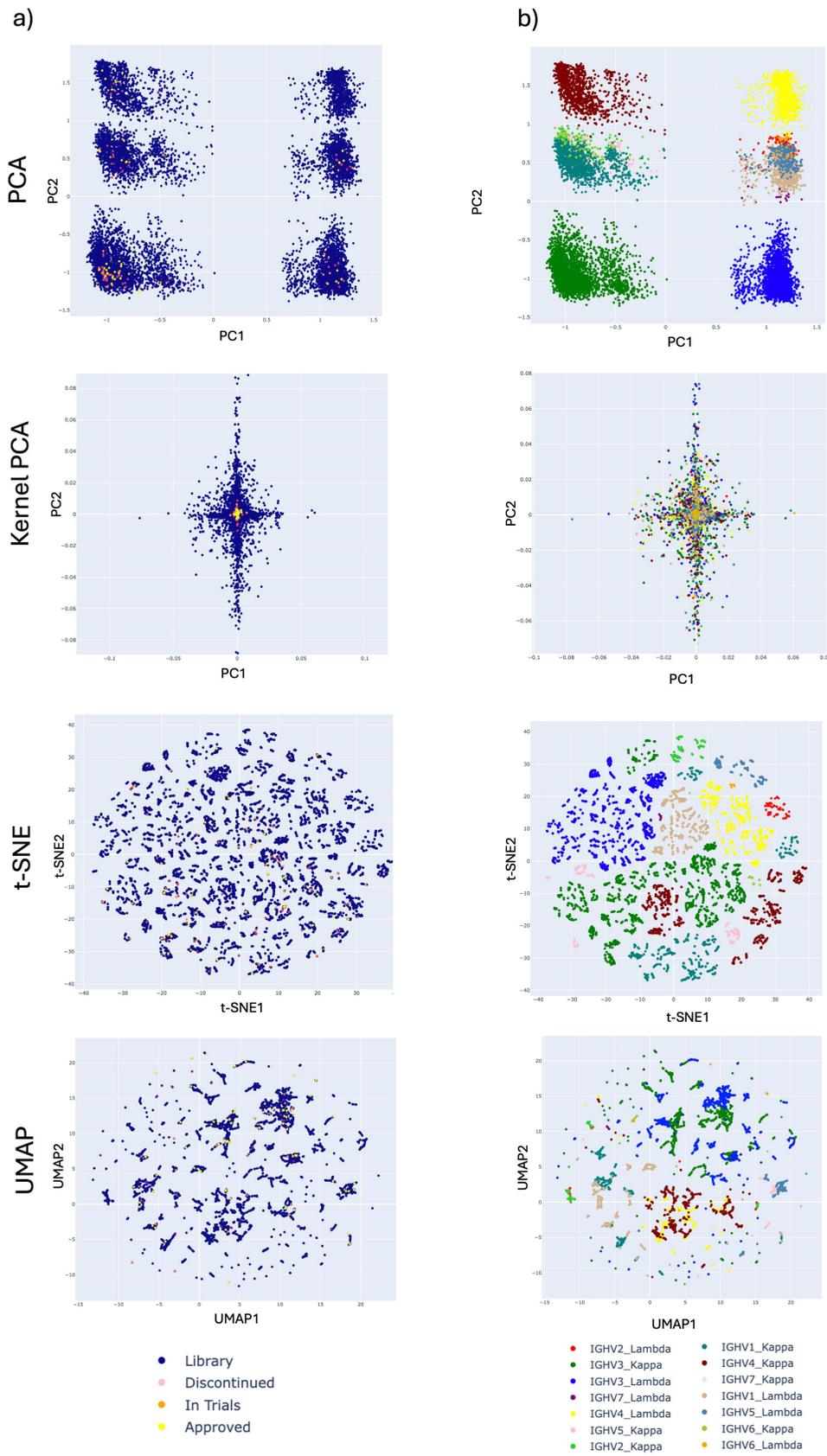


Figure 1. Scatter plots of unsupervised machine learning models trained on clinical ($n=144$) and library ($n=10,000$) paired antibody sequences encoded with the AntiBERTy protein language model. Plots are color coded by a) clinical stage or b) heavy chain V region germline gene and light chain type (λ or κ).

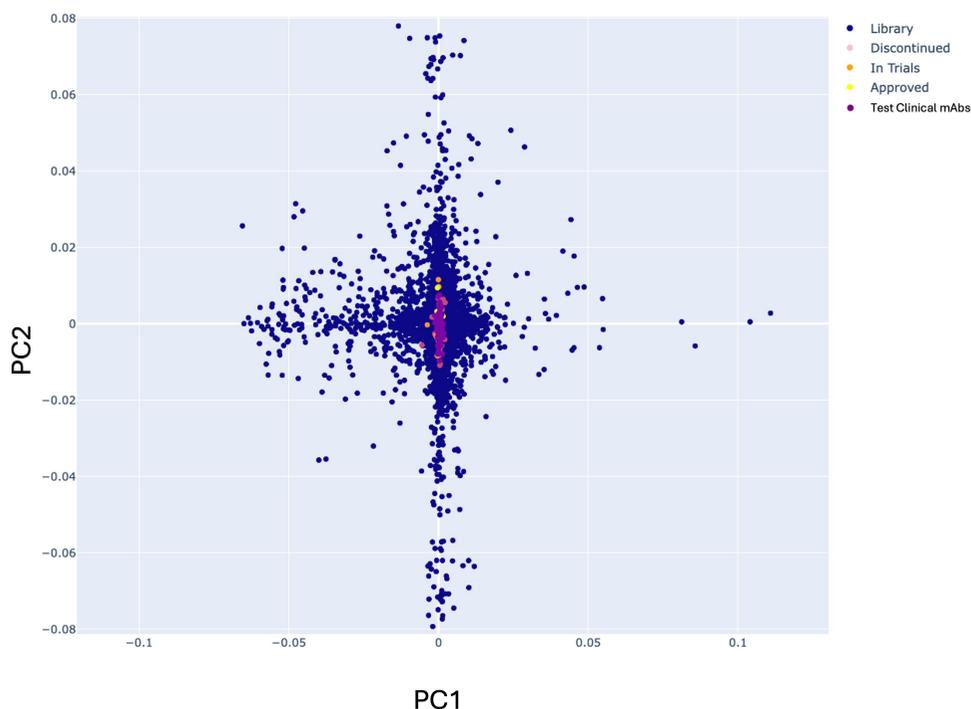


Figure 2. Scatter plot of kernel PCA (kernel='rbf', $\gamma = 500$) clinical mAbs trained on clinical ($n=144$), library ($n=10,000$) and a held-back test set of clinical ($n=203$) paired human antibody sequences encoded with the AntiBERTY language model.

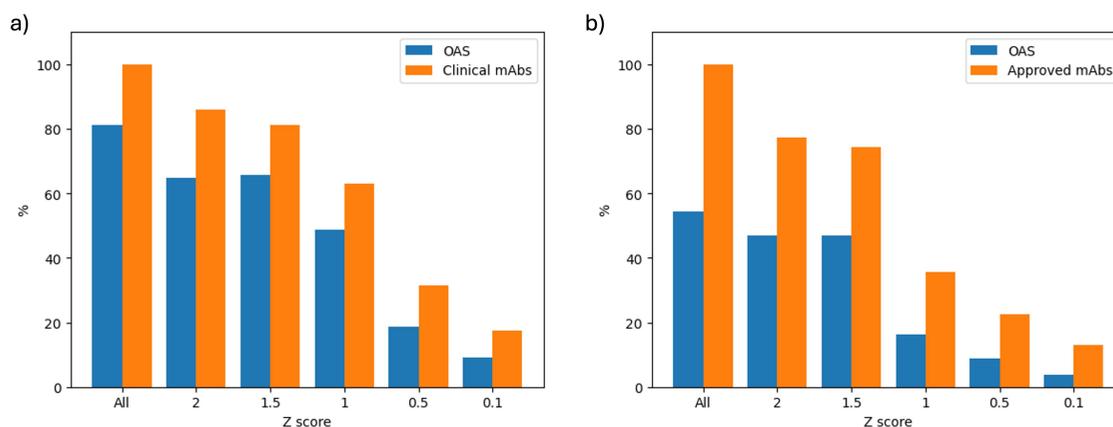


Figure 3. Percentages of OAS (library) and a) human clinical mAbs of any developmental stage, or b) only those with market approval, captured by the ellipse function drawn from the distribution of clinical mAbs. Z-scores denote how wide the distributions for the major axis of ellipse may be drawn with 'All' representing a Z-score selected such that all of the clinical (or approved, respectively) antibodies are captured.

Table 2. Means of sequence-calculated physicochemical properties for all market approved and discontinued mAbs (including human, humanized, chimeric and murine). p-values were calculated using a Mann-Whitney U test.

| Feature | Approved | Discontinued | p-value |
|--|-------------|--------------|---------|
| CDR-H3 Loop Length | 13.4 ± 4.25 | 10.7 ± 3.36 | 0.17 |
| $V_H \Delta G$ (kJ mol ⁻¹) | 7008 ± 3806 | 7592 ± 3424 | 0.35 |
| $V_L \Delta G$ (kJ mol ⁻¹) | 2411 ± 1351 | 2546 ± 2675 | 0.33 |
| Concatenated V_H/V_L (kJ mol ⁻¹) | 8434 ± 4855 | 1071 ± 4094 | 0.49 |
| Mean pI of V_H/V_L | 8.3 ± 1.18 | 7.9 ± 1.21 | 0.30 |
| Mean Minimum G score | -1.0 ± 1.22 | -0.8 ± 1.06 | 0.23 |

Details of the G score are given in Thullier *et al.*⁵⁶.

AbLang,⁴⁸ Sapiens⁴⁹) protein language models. The encodings of the paired V_H and V_L sequences were concatenated and treated as a single set of data points per antibody. The encoded antibody sequences were used to train a selection of 15

supervised machine learning classifiers (see the Methods and Supplementary File: Supplementary ML.pdf). Models were trained with ten-fold cross validation (CV) and model performance was evaluated using the mean Matthews' Correlation

Coefficient (MCC)⁵⁷ of the predictions of the test split of each fold. F-regression was used as a method of feature selection, by selecting the k features most correlated with the labels where k was set to [1,10,50,100,500,1000,2500,5000,10000].

Generally performance across all classifiers was good (Supplementary Tables S7–10), but overall the best performance was obtained for the AntiBERTy encodings, particularly when using F-regression for feature selection with k set to 2500. The top-performing models were the Linear Support Vector Machine classifier (LinearSVC; $\text{MCC} = 0.8 \pm 0.08$), Ridge Classifier ($\text{MCC} = 0.78 \pm 0.12$) and Logistic regression ($\text{MCC} = 0.80 \pm 0.1$) (See Figure 4). All methods were evaluated using a standard classification threshold of 0.5. The LinearSVC model was selected as the best model with a mean sensitivity of 0.86 ± 0.10 and specificity of 0.93 ± 0.05 across the 10 CV splits. In an attempt to improve the specificity further, this model was also assessed using a higher prediction threshold of

0.8. As expected, this resulted in a loss in sensitivity and an increase in specificity ($\text{Sn} = 0.57 \pm 0.17$, $\text{Sp} = 0.99 \pm 0.03$). This was accompanied by a decreased, but still respectable, MCC (0.64 ± 0.11). See Table 3 and Supplementary Figure S3, which shows confusion matrices for the raw outputs of this model at both probabilities.

The location of the majority of selected features in the V_H or V_L sequences was then identified. Intuitively, it was expected that CDR-H3 would contain a high proportion of features, but this was not the case. Instead, a region in Framework 3 of the V_L chain had a high concentration of selected features, indicating that this region is highly important in how all the models have learned to distinguish these groups (Figure 5).

However, when a ‘held back’ dataset of therapeutics which have been approved ($n=10$) and discontinued ($n=11$) since the original access of TheraSabDab (Supplementary Table S11), the predictive score was found to be $\text{MCC} = 0.14$ using

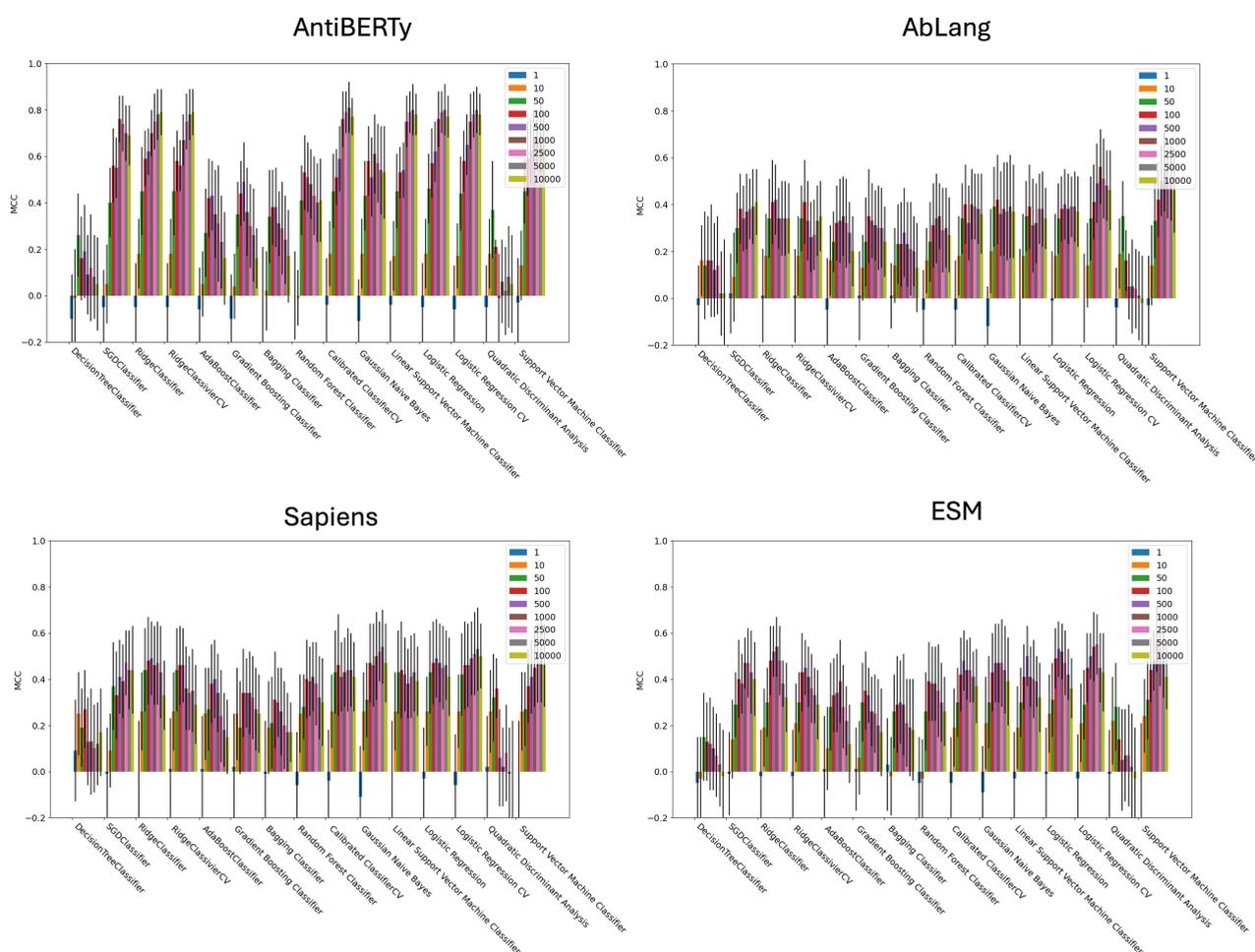


Figure 4. Matthews' correlation coefficient (MCC) and standard deviation from ten-fold cross validation of 15 binary machine learning predictors classifying approved ($n=115$) and discontinued ($n=150$) therapeutic antibodies and encoded using four protein language models.

Table 3. Summary performance of the LinearSVC supervised machine learning predictor for success in clinical trials.

| | Prediction Threshold | MCC | Performance Sensitivity | Specificity |
|------------------|----------------------|-----------------|-------------------------|-----------------|
| Cross-validation | 0.5 | 0.80 ± 0.08 | 0.86 ± 0.10 | 0.93 ± 0.05 |
| | 0.8 | 0.64 ± 0.11 | 0.57 ± 0.17 | 0.99 ± 0.03 |
| Independent | 0.5 | 0.14 | 0.50 | 0.64 |
| | 0.8 | 0.51 | 0.40 | 1.00 |

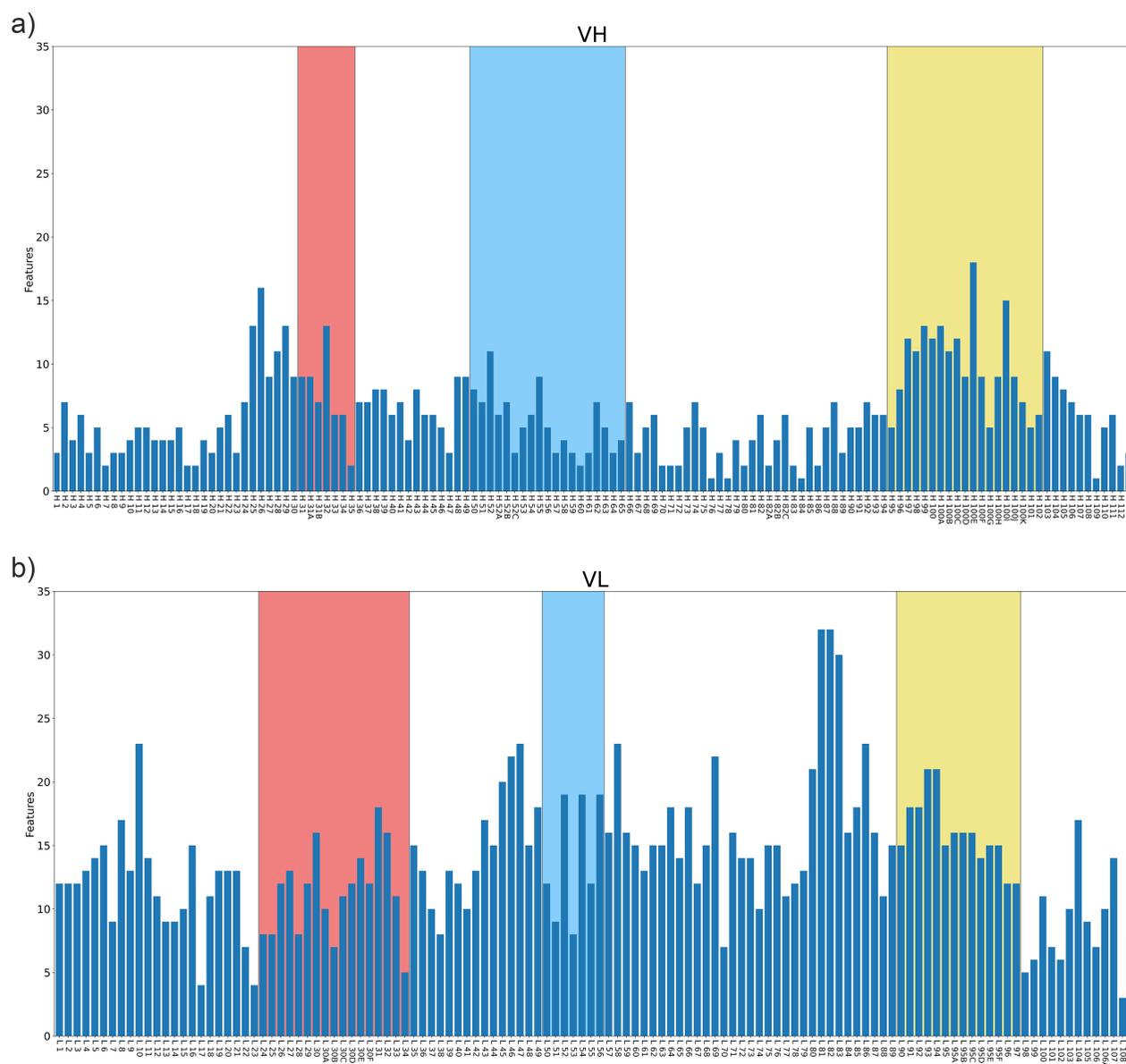


Figure 5. Locations of the top k features selected by F-regression from V_H and V_L chains of approved and discontinued mAbs encoded with the AntiBERTy language model⁵⁰ where $k = 2500$. CDR loops (Kabat definition) are highlighted in red (CDR1), Blue (CDR2) and Yellow (CDR3) with Chothia numbering.

a prediction threshold of 0.5, but this increased to $MCC = 0.51$ with the higher prediction threshold of 0.8 (Supplementary Figure S4a). A summary of results for the cross-validation and independent test sets is shown in Table 3. On the independent test set, the default prediction threshold results in a reasonable sensitivity and specificity as a result of numerous false positives; increasing the threshold to 0.8 improves the specificity accompanied by a small decrease in sensitivity, resulting in a much improved MCC. Supplementary Figure S4b provides confusion matrices, MCC, sensitivity and specificity at prediction thresholds between 0.5 and 0.9.

That such good performance was achieved using basic predictive models was surprising. Because approved and discontinued antibodies did not show a statistically significant difference in features, including isoelectric point, thermostability and CDR-H3 length (Table 2), other properties, such as immunogenicity²⁵ or the ability to access their targets, could be responsible for this ability to discriminate between these

groups.⁵⁰ It is also possible that subtle differences in V-region germline family pairing are involved, but Supplementary Figure S2 shows that these are broadly similar between the two groups.

Assembling the pipeline to optimize performance and offer additional triaging

The kernel PCA model ($\gamma = 500$) shown to separate library antibodies which are positioned close to clinical mAbs, and the LinearSVC model using 2500 features shown to separate approved and discontinued antibodies using the AntiBERTy language model encodings, were used to build a pipeline capable of selecting developable antibodies from an input library. The encoding only needs to take place once and can be carried forward between the two layers of models: unsupervised and supervised, respectively.

The complete pipeline attempts first to remove antibodies with obvious developability issues through physicochemical properties using Z-scores taken from the values of the approved antibody dataset ('Physicochemical Filtering' as given in Table 1, default $Z = 2$); this saves computational time in numbering and encoding antibodies with the language model, as well as producing a better quality output. Antibodies with features typical of clinical mAbs are then selected from the unsupervised clustering of encoded antibodies using the ellipse function ('Layer 1') and are then classified according to whether they are likely to pass clinical trials ('Layer 2'). A user may enter a library of human antibodies and obtain entries from those that are most likely to be successful. A schematic of the pipeline is shown in Figure 6 where stringency may be altered at each triaging step.

Testing on an example dataset demonstrates points of parameter tuning for optimized output

To illustrate the application of our pipeline, a library of 10,382 paired B-cell receptor (BCR) sequences taken from six healthy blood donors⁵⁹ was used as an example test dataset.

After physicochemical triaging, 'Layer 1' filtering is performed by performing PCA on the test data together with our 'library' antibodies from OAS and the previously used clinical dataset. The Z-score cutoffs are then calculated from the clinical dataset and the ellipse is generated and used to select antibodies from the test dataset.

Using decreasing Z-scores for the physicochemical triaging of the sequences reduced the number of antibodies entering 'Layer 1' (Table 4 and Supplementary Table S11). Similarly, decreasing the Z-score of the ellipse function in 'Layer 1' generally reduces the number of sequences taken forward to 'Layer 2' (Figure 7). However, since the clustering is performed and the ellipse is recalculated for each dataset, there is some variation and, in one case (Table 4, no physicochemical filtering, 'Layer 1', $Z = 1.0$), there is a small rise in the number of antibodies compared with $Z = 2.0$. Increasing the prediction threshold used in 'Layer 2' also reduces the final number of selected antibodies.

As a comparison for the quality of antibodies output by the model, we checked the TAP score³¹ for each antibody from the test BCR library. The TAP score is a developability score where an antibody with values for selected physicochemical properties that are seen within the clinical mAb dataset are given a perfect score of 0, and antibodies with increasing numbers of 'amber penalties' where the values are at the extremes of, or outside ('red flags'), the observed ranges are given negative scores. This is an indicator of developability, not whether an antibody is likely to be approved. It should be noted that, while the aims are broadly the same, this is a very different approach from our unsupervised machine learning: TAP relies solely on calculated or predicted physicochemical properties, while we use a subset of these properties only for a preliminary screen before using a clustering in high-dimensional space obtained from a protein language model.

The median TAP score for antibodies in the Test BCR library was 0, which means that more than half of these antibodies were predicted to have no developability warnings or

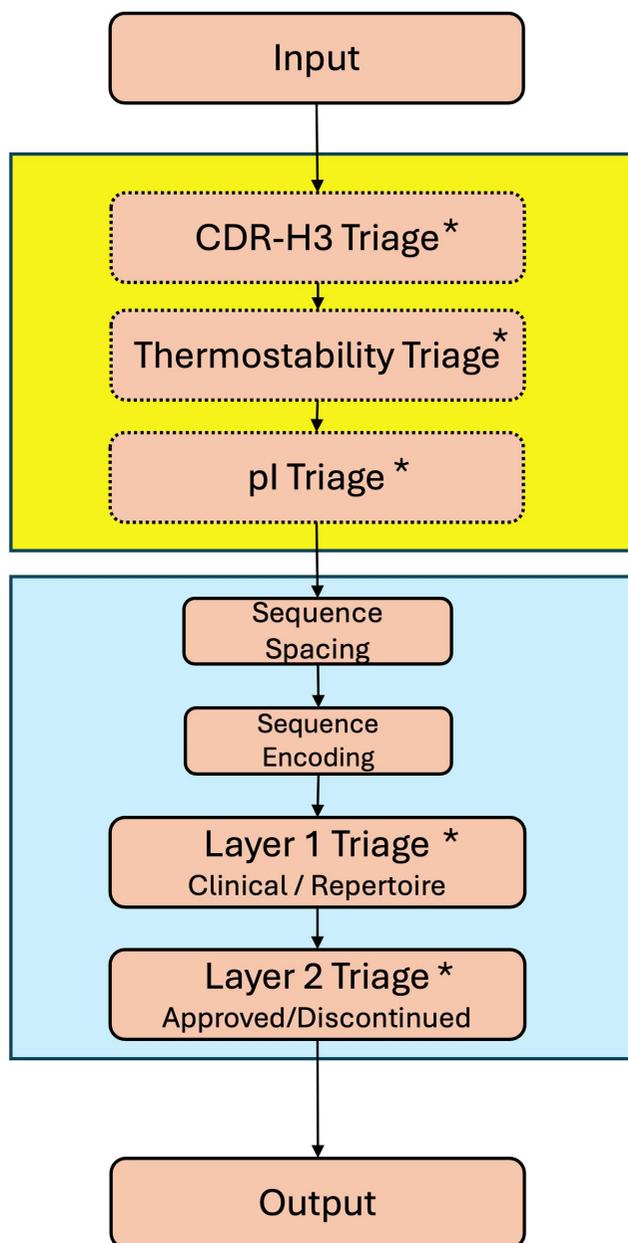


Figure 6. Schematic of the antibody triaging pipeline. The yellow box indicates optional physicochemical feature triaging steps calculating CDR-H3 length using AbNum⁵⁸ Thermostability (ΔG of unfolding) is calculated using the Oobatake Method⁴⁵ and pI using the IPC method.⁴⁶ the blue box indicates machine learning elements including spacing and encoding, as well as 'Layer 1' triage which is based on the kernel PCA model for separating antibodies with similar properties to clinical mAbs from the repertoire. The selection of antibodies to take forward is made using the ellipse function. 'Layer 2' is the supervised LinearSVC model trained to distinguish approved and discontinued clinical mAbs. '*' indicates stages where stringency can be adjusted using Z-score thresholds, or the prediction threshold in the case of 'Layer 2'.

red flags. However, the minimum TAP score observed from the library was -110 , indicating there are antibodies in the library with many developability warnings or red flags. From the data in Table 4, it is clear that setting the physicochemical property filtering (PCF) in our approach to a more stringent Z-score (e.g., $Z = 0.5$) had the major effect in removing antibodies with the most negative TAP scores from the output. Indeed with no PCF, neither the 'Layer 1' nor 'Layer 2' filtering removed the antibodies with $TAP = -110$. Similarly, as the

Table 4. Number of antibodies from the test BCR library output from the triaging pipeline given different parameters of physicochemical filtering (PCF) and 'Layer 1' thresholds. For comparison, the minimum and mean TAP scores are provided, together with the percentage of negative TAP scores, after 'Layer 1' and 'Layer 2' shown separated by a '/'.

| PCF Z-score | | Layer 1 Filtering Z-Score | | | |
|-------------|------------------|---------------------------|---------------|---------------|---------------|
| | | None | 2.0 | 1.0 | 0.5 |
| None | PCF Only | 10492 | – | – | – |
| | Layer 1 | 9875 | 8165 | 8186 | 6107 |
| | Layer 2 | 3587 | 2981 | 2978 | 2232 |
| | Min TAP Score | –110/–110 | –110/–110 | –110/–110 | –110/–110 |
| | Mean TAP Score | –18.58/–18.86 | –18.52/–18.91 | –18.53/–18.97 | –18.28/–18.83 |
| 2.0 | % TAP Scores < 0 | 40.3/50.4 | 48.0/50.4 | 30.4/50.2 | 22.2/50.7 |
| | PCF Only | 8045 | – | – | – |
| | Layer 1 | 7508 | 7333 | 5855 | 3753 |
| | Layer 2 | 2571 | 2514 | 1981 | 1272 |
| | Min TAP Score | –110/–110 | –110/–110 | –110/–110 | –110/–90 |
| 1.0 | Mean TAP Score | –18.07/–18.40 | –18.05/–18.47 | –18.12/–18.50 | –18.11/–17.58 |
| | % TAP Scores < 0 | 44.2/47.1 | 44.1/47.2 | 43.9/46.8 | 19.7/47.4 |
| | PCF Only | 2740 | – | – | – |
| | Layer 1 | 2359 | 2329 | 2056 | 1086 |
| | Layer 2 | 808 | 797 | 705 | 361 |
| 0.5 | Min TAP Score | –90/–40 | –90/–40 | –90/–40 | –90/–40 |
| | Mean TAP Score | –16.77/–57 | –16.78/–16.43 | –16.83/–16.33 | –16.99/–16.23 |
| | % TAP Scores < 0 | 31.7/35.3 | 31.7/35.5 | 32.3/35.6 | 31.8/36.0 |
| | PCF Only | 386 | – | – | – |
| | Layer 1 | 308 | 231 | 157 | 39 |
| | Layer 2 | 113 | 80 | 57 | 14 |
| | Min TAP Score | –40/–40 | –40/–40 | –30/–30 | –20/–20 |
| | Mean TAP Score | –15.68/–15.0 | –15.25/–5.38 | –15.33/–15.00 | –11.0/–12.5 |
| | % TAP Scores < 0 | 25.3/31.9 | 21.5/32.5 | 28.7/38.6 | 35.6/38.6 |

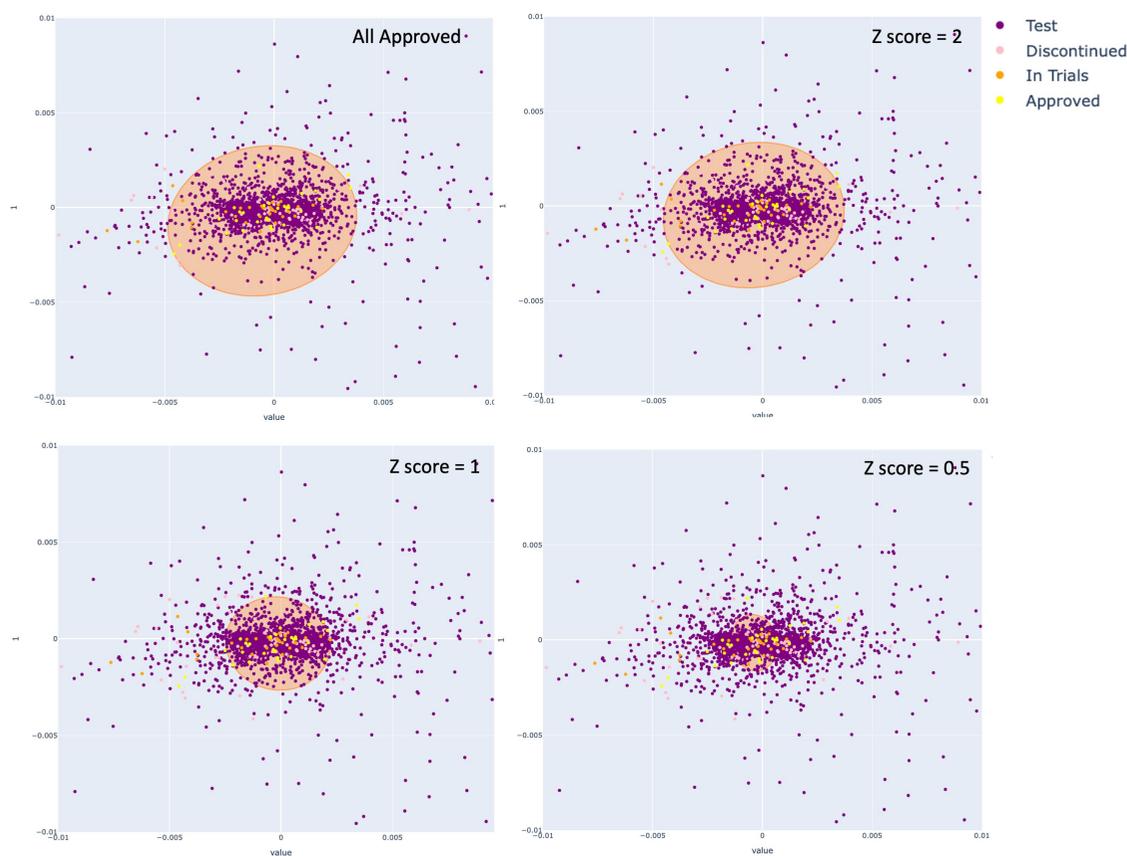


Figure 7. Scatter plots of clinical ($n=144$) and library ($n=2740$) paired antibody sequences encoded with the AntiBERTy protein language model and that have undergone dimensionality reduction using kernel Principal component analysis with a radial basis kernel function ($\gamma = 500$). Different Z-scores of the distribution of clinical antibodies along PC1 are used as the extremes of the major axis to draw the ellipse function.

'Layer 1' stringency was increased, there was very little effect on the minimum, or the mean, TAP score. This is, perhaps, not surprising since the physicochemical properties on which this preliminary filtering is performed are somewhat similar to those exploited by the TAP score. However, the number of negative TAP scores does decrease as the 'Layer 1' filtering becomes more stringent (Table 4).

Again, because of the recalculation of the Z-scores and ellipse, there is one case in which the mean TAP score does not steadily progress closer to zero as the 'Layer 1' stringency is increased (Table 4, physicochemical filtering, Z-score = 0.5).

It is also interesting to observe that, comparing the output of 'Layer 1' and 'Layer 2', the minimum and mean TAP scores improve. Given that 'Layer 2' is predicting clinical success rather than developability, there is no reason to expect that this would be the case. Indeed, the percentage of antibodies with negative TAP scores retained after 'Layer 2' is larger than that after 'Layer 1', indicating that 'Layer 2' filtering is indeed detecting something different from developability.

Discussion

We have demonstrated the ability to triage library antibodies to find those with properties similar to currently available therapeutic mAbs. This was achieved through a combination of preliminary filtering using physicochemical properties (to remove clearly outlying mAbs), with unsupervised and supervised machine learning. This demonstrates a useful tool in monoclonal antibody therapeutic discovery that may be applied to new and preexisting paired human antibody libraries to identify possible clinical candidates with potential to pass clinical trials in order to avoid expensive late-stage failures. Parameters of the pipeline at each step may be adjusted such that increased or reduced stringency filtering can produce a smaller (but more likely to be successful) or larger selection of antibodies. This pipeline can be used to identify antibodies with properties of therapeutic mAbs from large libraries,^{60,61} to screen antibodies from transgenic animals following immunizations,⁶² or from human patients recovering from a condition of interest.⁵ Using the pipeline in these contexts reduces the experimental work in finding an antibody which has properties suitable for use in the clinic.²⁷

The 'Layer 1' triage creates a 2D projection of the N -dimensional space resulting from the protein language model encoding. We found that all the clinical stage antibodies clustered in this space, and that the held-back set also clustered with the original set of clinical antibodies. Consequently, this must represent a region of the N -dimensional space that shares some properties and therefore other antibodies that are found in this region must have similar features. Nonetheless, it is perfectly possible that some antibodies that do not have suitable developability properties for use in the clinic may also fall within the bounds that we set around the cluster; it is also likely that some future clinically successful antibodies may not cluster well with the existing clinical set. The purpose of this step is to identify candidates with the highest chance of having good developability characteristics.

The pipeline allows for further optional triaging to be added at any point to give additional layers of stringency. The

advantage of using these steps is a vastly reduced computation time. On an A5000 GPU used in this study, the AntiBERTy encoding takes 0.06 seconds per V_H and V_L pair, making it suitable for the high-throughput analysis of libraries (compared with the approximately 30 seconds required per antibody for the TAP Score web server³¹). While the protein language models may be doing so implicitly, using additional features also opens up the possibility of using other explicit features, including screening for immunogenicity⁵⁵ and known sequence liabilities such as post-translation modification sites^{60,63} and hydrophobic patches.²⁴

Direct comparisons of the performance of our method with the TAP score are not really possible. TAP relies on the distribution of a number of calculated and predicted physicochemical properties, some of which rely on a (predicted) structure of the antibody. While predicted properties can be compared with available experimental data, we only use physicochemical properties (calculated solely from sequence) as a preliminary screen to remove obvious outliers. The machine learning stages are based on a protein language model encoding that projects information (including implied structural features) into a high-dimensional space, which is then reduced to a 2D space in which clinical antibodies are seen to cluster. Consequently, we do not directly predict properties related to developability and comparisons with published experimental data are not possible.

It is also worth noting that our approach is not simply suggesting that if the sequences are more similar to those of clinical-stage antibodies, they should have better developability. If that were the case, we could just have used a BLAST search. Rather, we exploit the encoding from the antibody-specific protein language model, AntiBERTy and it is well established that protein language model encoding of sequences relates a sequence to structural and lineage information on which it has been trained and thus captures other key information. These encodings are highly sensitive and can even predict the effect of single amino acid changes.⁶⁴ We have identified a 2D projection of the AntiBERTy encoding that clusters the clinical-stage antibodies and consequently we are looking at similarity of the protein language model encodings in those two principal components rather than sequence similarity *per se*.

Because human clinical antibodies clustered so closely in the kernel PCA, they must have similar features which have been encoded and recognized by the language model. The fact that the clinical antibodies cluster near the origin suggests that they are developable largely because they are 'ordinary' antibodies which innately satisfy the required conditions. As stated above, it is entirely possible that antibodies with very different properties could have therapeutic potential, but ultimately these would be higher-risk and consequently it is generally better to allow false negatives than false positives.

It is interesting to note that, using the TAP score, more than 25% of fully human clinical-stage antibodies exhibit negative TAP scores. Our approach clusters all the clinical antibodies and the ellipse function (with default parameters) will capture all of these, including those that have TAP red flags. In other words, if input antibodies are found to be located within the same region of the

projection of the high-dimensional encoding, they are likely to have *sufficiently good* developability. We calculated the TAP scores for the 133 human clinical-stage antibodies and found that, while ~74% have a TAP score of zero (indicating no developability issues), the remainder have negative TAP scores as low as -30 (see Supplementary Figure S5).

While the unsupervised model ('Layer 1') groups together both approved and discontinued clinical antibodies, when these two groups are studied in a supervised context ('Layer 2'), it is possible to recognize differences between them. Even though the dataset is small, the 10-fold cross-validation and held-back dataset demonstrate that there are features that are important for successfully completing clinical trials and, in future, larger datasets would allow us to have increased confidence in these predictions. The results also show that the light-chain Framework 3 seems to have a large contribution to these features. It is worth noting that this predictor is not identifying something trivial in the sequences. Supplementary Figure S2 shows that the distribution of germline light/heavy pairs is very similar in the antibodies that succeed or fail in clinical trials. As an additional control to look for simple sequence features, we also took a very simple approach of predicting that all antibodies with kappa light chains would succeed while all with lambda light chains would fail. As expected, this showed an MCC of 0.01 indicating no predictive power (data not shown).

A drug may be discontinued from trials for efficacy reasons relating to bioavailability or binding to the target, safety reasons relating to the antigen or antibody (including immunogenicity) as well as business or marketing reasons (including the existence of other good drugs).^{1,9} As discussed above, since we showed that there are no statistical differences between the approved and discontinued groups for thermostability, pI or CDR-H3 properties, it is possible that the model is selecting features related to immunogenicity, or V_H/V_L germline gene pairing which may be related to stability.⁶⁵ The latter option could then be related to biases seen in the approved and discontinued datasets perpetuated by the lead candidate selection processes, although the approved and discontinued dataset have similar proportions of V_H/V_L germline gene pairing (Supplementary Figure S2), indicating other factors are being recognized in this region which are related to clinical trial success.

To summarize, we have developed a tool with a goal similar to methods such as TAP,³¹ but that works in a different way, exploiting 'big data' for protein language model encodings and exploring a large sample of the human repertoire, together with artificial intelligence. This work has demonstrated the ability to triage a library of antibodies to identify those with features similar to approved mAb therapeutics (and rejecting those that are very different) using language model encoding and applying them to both unsupervised and supervised machine learning. Furthermore, we demonstrate the ability to fine-tune the output in terms of quality by adjusting the thresholds of the models used to obtain the output. These tools aim to make use of previously curated and future antibody

datasets to triage large datasets, enabling faster and cheaper identification of potential lead candidates.

Methods

Data collection

Human clinical-stage mAbs

Paired V_H and V_L sequences of therapeutic monoclonal antibodies ($n=801$) were downloaded from the October 2021 release of TheraSabDab.³ Therapeutics marked as 'Whole mAb' were selected and identified as being fully human using the '-umab' suffix excluding instances of '-zumab' (humanized). Each therapeutic was checked for its source using the literature. This resulted in a dataset of 143 antibodies: approved mAbs ($n=31$); discontinued mAbs ($n=77$) and in trials ($n=35$) (Supplementary Table S1). A further independent test dataset of human-derived clinical mAbs was acquired ($n=203$) using the 2016 naming convention in which the source infix was removed from the name and the 2022 naming convention using '-tug' for unmodified whole immunoglobulins and '-bart' for whole immunoglobulins with engineered amino acid changes in the constant domains⁵⁴ (Supplementary Table S4).

Library antibodies from OAS

The Observed Antibody Space database¹⁸ was accessed in January 2022 and 34 libraries were downloaded totaling 88,274 paired sequences. A total of 10000 antibodies were selected randomly in order to create a training set for unsupervised learning (Supplementary Table S2).

Approved and discontinued mAbs

Clinical mAbs were obtained from the October 2021 release of the TheraSabDab database.³ The V_H and V_L sequences of 115 approved antibody drugs and 156 discontinued drugs were collected. Seven drugs were excluded from the discontinued dataset as they were found to be discontinued for reasons not related to efficacy or safety. Edrecolomab was also moved from the approved dataset and the discontinued dataset because it was later withdrawn for efficacy reasons.⁶⁶ The result of this is a dataset of 115 approved and 150 discontinued antibodies (Supplementary Table S5). Excluded sequences and reasons for their exclusion are found in Supplementary Table S12. A held-back dataset of 21 therapeutics was taken from TheraSabDab accessed in October 2023 and not included in the original dataset (Supplementary Table S11).

Test BCR dataset

The Test B cell receptor (BCR) sequence dataset⁶⁷ used to demonstrate the pipeline was downloaded from dx.doi.org/10.5281/zenodo.5146019. This dataset was obtained from six healthy blood donors whose B cells were isolated and sorted via fluorescence-activated cell sorting by developmental stage. Transcripts from each individual cell were bar-coded making V_H/V_L pairing possible. Antibody V_H and V_L pairs were taken from B cells which shared the same bar-code where both an IGH and IG λ or IG κ chain was present. In cases where both

IG λ and IG κ chains were present, the chain with the highest count number was taken as the V_L chain pair. No filtering based on the type or BCR developmental stage was performed. Individual amino acid sequences for frameworks and complementarity-determining region (CDR) loops were concatenated to give the full antibody variable domain sequence. In total, 10,382 paired antibodies were extracted.

Encoding H and L sequences with antibody language models

V_H and V_L sequences were numbered according to the Chothia scheme⁶⁸ using AbNum⁵⁸ (www.bioinf.org.uk/abs/abnum/), where missing residues in the numbering scheme sequence were padded with characters dependent on which protein language model was being used, to align all sequences making V_H sequences 132 residues long and V_L sequences 122 residues long. Details of sequence encodings can be found in Table 5.

Supervised machine learning

Supervised learning was performed with SciKitLearn using 15 classifiers⁶⁹ given in Table 6. Descriptions of each classifier used and details can be found in Supplementary File: Supplementary_ML.pdf.

F-regression is a method of feature reduction where the k most informative features are kept as input to the model. This is done by calculating the cross-correlation of each data point and the label for all features, which is converted to an F-score, then to a p-value and ranked.⁷⁰ F-regression was implemented through the module `sklearn.feature_selection.SelectKBest` using the Python module `sklearn.feature_selection.f_regression` as the score

Table 5. Details of language model encodings.

| Language Model | Features (VH+VL) | Padding Character | Reference |
|-------------------------|------------------|-------------------|-----------|
| AntiBERTy | 130,048 | ' ' | 50 |
| AbLang | 195,072 | '*' | 18 |
| Sapiens | 152,560 | '*' | 49 |
| ESM (esm2 t6 8 M UR50D) | 82,560 | 'X' | 44 |

function and variable numbers for k were substituted [1,10,50,100,500,1000,2500,5000,10000].

Once the F-regression was implemented on the encoded dataset, it was then split into training and test sets using `sklearn.model_selection.train_test_split` where training portions were used to train the models using ten-fold cross-validation.

Model performance was measured using the Matthews' Correlation Coefficient (MCC),⁵⁷ which gives a score between -1 (perfect inverse prediction) and 1 (perfect prediction), with 0 being random chance. Mean MCC and standard deviation for prediction performance over the ten folds were reported.

Unsupervised machine learning

PCA was used as a method of dimensionality reduction and implemented through `sklearn.decomposition.PCA`. Non-linear (kernel) PCA⁵¹ was implemented through `sklearn.decomposition.KernelPCA` using kernel functions 'rbf', 'cosine' and 'poly' and two principal components. At first the coefficient of the kernel (γ) was set to the default value of $1/k$ where k is the number of features. Once rbf was selected as the most suitable method, differing values for γ were tested [10, 50,100, 500, 1000]. t-distributed Stochastic Neighbor Embedding (t-SNE)⁵² was implemented through `sklearn.manifold.TSNE` with two components where the learning rate was set to 10, and the perplexity set to 1000. Uniform manifold approximation and projection (UMAP)⁵³ was implemented through `sklearn.manifold.UMAP` with the learning rate set to 1 and the nearest neighbors set to 100.

Ellipse function

The ellipse function takes in the points of the two extremes on the major axis ($x1, y1$) and ($x2, y2$), as well as a value for h (the height of the minor axis). The major axis is taken as the principal component where clinical mAbs have the largest distribution, and the selected points are given as the points on the distribution closest to a given Z-score in that distribution. The value of h is given as the distance between the two equivalent points on the minor axis. The method for producing the ellipse works as follows:

Table 6. Supervised machine learning classifiers used in classifying approved and discontinued antibodies.

| Classifier | Acronym | Implementation |
|--|------------|--|
| Decision Tree | | <code>sklearn.tree.DecisionTreeClassifier</code> |
| Stochastic Gradient Descent Classifier | SGDC | <code>sklearn.linear_model.SGDClassifier</code> |
| Ridge Classifier | | <code>sklearn.linear_model.RidgeClassifier</code> |
| Ridge Classifier CV | | <code>sklearn.linear_model.RidgeClassifierCV</code> |
| AdaBoost Classifier | | <code>sklearn.ensemble.GradientBoostingClassifier</code> |
| Gradient Boost Classifier | | <code>sklearn.ensemble.GradientBoostingClassifier</code> |
| Bagging Classifier | | <code>sklearn.ensemble.BaggingClassifier</code> |
| Random Forest Classifier | | <code>sklearn.ensemble.RandomForestClassifier</code> |
| Calibrated Classifier | | <code>sklearn.calibration.CalibratedClassifier</code> |
| Gaussian Naive Bayes Classifier | GaussianNB | <code>sklearn.naive_bayes.GaussianNB</code> |
| Support Vector Machine | SVC | <code>sklearn.svm.SVC</code> |
| Linear Support Vector Machine Classifier | LinearSVC | <code>sklearn.svm.LinearSVC</code> |
| Logistic Regression Classifier | | <code>sklearn.linear_model.LogisticRegression</code> |
| Logistic Regression CV Classifier | | <code>sklearn.linear_model.LogisticRegressionCV</code> |
| Quadratic Discriminant Analysis Classifier | QDA | <code>sklearn.discriminant_analysis.QuadraticDiscriminantAnalysis</code> |

- Calculate the major and minor radii of the ellipse (a and b , respectively). The major radius is calculated from the two given points (Equation 1) and the minor radius is calculated as half the value given for h , where Δx is the difference in x values and Δy is the difference in y values between the two extreme points on the major axis.

$$a = \frac{\sqrt{\Delta x^2 + \Delta y^2}}{2}, b = \frac{h}{2} \quad (1)$$

- Use the parametric equation of an ellipse to generate the ellipse over 100 equally spaced points between 0 and 2π assuming it is centered at the origin (Equation 2). For a given point on the ellipse:

$$x = a \cos(\theta), y = b \sin(\theta) \quad (2)$$

where a is the major axis radius, b is the minor axis radius and θ is a given angle between 0 and 2π .

- Calculate the angle between given points to obtain an angle of rotation using the Python Numpy arctan2 function⁷¹ for Δy and Δx .
- Calculate a rotation matrix (R) based on the angle of rotation:

$$R = [[\cos(\theta), -\sin(\theta)], [\sin(\theta), \cos(\theta)]] \quad (3)$$

where θ is the angle of rotation.

- Apply the rotation matrix R, to the ellipse.
- Calculate the midpoint of the two given points:

$$x = \frac{x_1 + x_2}{2}, y = \frac{y_1 + y_2}{2} \quad (4)$$

- Translate the ellipse to the midpoint.
- For each point, check if its x and y coordinates are inside the ellipse using the Polygon function from the Python 'Shapely' module.

Calculating physicochemical properties

Physicochemical properties were calculated as described below and compared between groups using the two-tailed unpaired Mann-Whitney U-test.⁷²

Identifying CDR-H3 loops

The CDR loop three of the V_H domain (CDR-H3) has frequently been observed to have the largest contribution to antibody binding affinity because it is the most diverse region between sequences, overlapping the Variable, Diversity and Junction gene segments.^{47,73} CDR-H3 regions were identified using the AbNum software⁵⁸ and applying the Kabat/Chothia/Martin definition (H95-H102). Sequences with more than two cysteine residues were excluded as additional cysteines are a known risk factor for aggregation.⁷⁴

Thermostability

Gibbs Free Energy (ΔG) of unfolding was predicted for each antibody sequence using the Oobatake method⁴⁵ with

experimental values of ΔH and ΔS taken from the original paper. mAbs with negative ΔG of unfolding values were considered unstable and associated with poor developability. This was calculated for the V_H and V_L chains, as well as for both chains concatenated together using the 'ssbio' Python module.⁷⁵

Isoelectric point

The method of calculating Isoelectric Point (pI) was that used in the IPC software⁴⁶ which uses experimentally obtained peptide pKa values from the EMBOSS database⁷⁶ substituted into a rearranged Henderson-Hasselbach equation. The equations are iterated using different pH values, starting at 6.5, and the results of the termini and each of the charged residues are summed together. If the sum is 0 ± 0.01 , the isoelectric point is reached. Otherwise, the iteration continues to increase the pH if the summed net charge was positive or to decrease the pH if it was negative.

Immunogenicity (humanness)

The G score⁵⁶ is a measure of antibody humanness based on similarity to germline families and a predictor of immunogenicity. This metric was calculated using the online tool www.bioinf.org.uk/abs/gscore/ for V_H and V_L independently. The minimum score of these chains for each antibody was taken and the mean for each of these sets of minima is presented.

V-region germline gene identification

V-region Germline genes were identified using the in-house 'Assign Germline' software (AGL; github.com/AndrewCRMartin/agl/). Where more than one germline gene has the same (highest) sequence identity, AGL selects a gene using the logic that the germline family with the lowest family number was likely to have been discovered first and therefore likely to be more numerous. The same logic is applied to allelic variants and proximal genes are favored over distal genes, ensuring that gene names are consistent.

TAP scores

TAP scores were developed by Raybould *et al.*³¹ to compare an antibody with the clinical dataset using metrics related to developability, assigning 'amber penalties' to antibodies that fall in the top and bottom 5% of the observed distribution, and 'red flags' to antibodies that fall outside the distribution. TAP scores were calculated for 10,382 paired V_H and V_L nucleotide sequences from the Test BCR dataset in batches of 500 using the IGX platform igx.bio/ in August 2023 using the default penalty set. Details of statistics measured and penalties assigned can be found in Raybould *et al.*³¹

Acknowledgments

The test library of paired B cell receptor (BCR) sequences was supplied by Franca Fraternali and Joseph Ng (Institute of Structural and Molecular Biology, UCL). TAP scores were calculated using the IGX platform developed by ENPICOM (NL) with academic licence.

Disclosure statement

No potential conflict of interest was reported by the author(s).

Funding

This work was supported by the BBSRC LIDo program under Grant [BB/T008709/1].

ORCID

Andrew C.R. Martin  <http://orcid.org/0000-0002-2835-2572>

References

- Khetan R, Curtis R, Deane CM, Hadsund JT, Kar U, Krawczyk K, Kuroda D, Robinson SA, Sormanni P, Tsumoto K, et al. Current advances in biopharmaceutical informatics: guidelines, impact and challenges in the computational developability assessment of antibody therapeutics. *mAbs*. 2022;14(1):2020082. doi:10.1080/19420862.2021.2020082.
- Alfaleh MA, Alsaab HO, Mahmoud AB, Alkayyal AA, Jones ML, Mahler SM, Hashem AM. Phage display derived monoclonal antibodies: from bench to bedside. *Front Immunol*. 2020;11:1986. doi:10.3389/fimmu.2020.01986.
- Raybould MIJ, Marks C, Lewis AP, Shi J, Bujotzek A, Taddese B, Deane CM. Thera-SabDab: the therapeutic structural antibody database. *Nucleic Acids Res*. 2020;48(D1):D383–D388. doi:10.1093/nar/gkz827.
- Taylor PC, Adams AC, Hufford MM, de la Torre I, Winthrop K, Gottlieb RL. Neutralizing monoclonal antibodies for treatment of COVID-19. *Nat Rev Immunol*. 2021;21(6):382–393. doi:10.1038/s41577-021-00542-x.
- Keam SJ. Tixagevimab + cilgavimab: first approval. *Drugs*. 2022;82(9):1001–1010. doi:10.1007/s40265-022-01731-1.
- Kaplon H, Chenoweth A, Crescioli S, Reichert JM. Antibodies to watch in 2022. *mAbs*. 2022;14(1):2014296. doi:10.1080/19420862.2021.2014296.
- Kaplon H, Crescioli S, Chenoweth A, Visweswarajah J, Reichert JM. Antibodies to watch in 2023. *mAbs*. 2023;15(1):2153410. doi:10.1080/19420862.2022.2153410.
- Troisi M, Marini E, Abbiento V, Stazzoni S, Andreano E, Rappuoli R. A new dawn for monoclonal antibodies against antimicrobial resistant bacteria. *Front Microbiol*. 2022;13:1080059. doi:10.3389/fmicb.2022.1080059.
- Sun A, Benet LZ. Late-stage failures of monoclonal antibody drugs: a retrospective case study analysis. *Pharmacology*. 2020;105(3–4):145–163. doi:10.1159/000505379.
- Robinson WH. Sequencing the functional antibody repertoire—diagnostic and therapeutic discovery. *Nat Rev Rheumatol*. 2015;11(3):171–182. doi:10.1038/nrrheum.2014.220.
- Rajan S, Kierny MR, Mercer A, Wu J, Tovchigrechko A, Wu H, Dall'Acqua WF, Xiao X, Chowdhury PS. Recombinant human B cell repertoires enable screening for rare, specific, and natively paired antibodies. *Commun Biol*. 2018;1(1):5. doi:10.1038/s42003-017-0006-2.
- Jaffe DB, Shahi P, Adams BA, Chrisman AM, Finnegan PM, Raman N, Royall AE, Tsai F, Vollbrecht T, Reyes DS, et al. Functional antibodies exhibit light chain coherence. *Nature*. 2022;611(7935):352–357. doi:10.1038/s41586-022-05371-z.
- Irac SE, Soon MSF, Borcharding N, Tuong ZK. Single-cell immune repertoire analysis. *Nat Methods*. 2024;21(5):777–792. doi:10.1038/s41592-024-02243-4.
- Martin ACR. Accessing the Kabat antibody sequence database by computer. *Proteins Struct Funct Bioinf*. 1996;25(1):130–133. doi:10.1002/(SICI)1097-0134(199605)25:1<130::AID-PROT11>3.0.CO;2-L.
- Lefranc M-P. IMGT, the International ImMunoGeneTics information system. *Cold Spring Harbor Protocol*. 2011;2011(6):595–603. doi:10.1101/pdb.top115.
- Dunbar J, Krawczyk K, Leem J, Baker T, Fuchs A, Georges G, Shi J, Deane CM. SabDab: the structural antibody database. *Nucleic Acids Res*. 2013;42(Database issue):D1140–D1146. doi:10.1093/nar/gkt1043.
- Ferdous S, Martin ACR. AbDb: antibody structure database — a database of pdb-derived antibody structures. *Database*. 2018;2018:bay040. doi:10.1093/database/bay040.
- Olsen TH, Boyles F, Deane CM. Observed antibody space: a diverse database of cleaned, annotated, and translated unpaired and paired antibody sequences. *Protein Sci*. 2022;31(1):141–146. doi:10.1002/pro.4205.
- Guo Y, Chen K, Kwong PD, Shapiro L, Sheng Z. cAb-Rep: a database of curated antibody repertoires for exploring antibody diversity and predicting antibody prevalence. *Front Immunol*. 2019;10:2365. doi:10.3389/fimmu.2019.02365.
- Margreitter C, Lu H-C, Townsend C, Stewart A, Dunn-Walters DK, Fraternali F. Brepertoire: a user-friendly web server for analysing antibody repertoire data. *Nucleic Acids Res*. 2018;46(W1):W264–W270. doi:10.1093/nar/gky276.
- Fernández-Quintero ML, Loeffler JR, Kraml J, Kahler U, Kamenik AS, Liedl KR. Characterizing the diversity of the CDR-H3 loop conformational ensembles in relationship to antibody binding properties. *Front Immunol*. 2019;9:3065. doi:10.3389/fimmu.2018.03065.
- Wolf Perez A-M, Lorenzen N, Vendruscolo M, Sormanni P. Assessment of therapeutic antibody developability by combinations of in vitro and in silico methods. *Met Mol Biol*. 2022;2313:57–113.
- Obrezanova O, Arnell A, de la Cuesta RG, Berthelot ME, Gallagher TRA, Zurdo J, Stallwood Y. Aggregation risk prediction for antibodies and its application to biotherapeutic development. *mAbs*. 2015;7(2):352–363. doi:10.1080/19420862.2015.1007828.
- Waibl F, Fernández-Quintero ML, Wedl FS, Kettenberger H, Georges G, Liedl KR. Comparison of hydrophobicity scales for predicting biophysical properties of antibodies. *Front Mol Biosci*. 2022;9:960194. doi:10.3389/fmolb.2022.960194.
- Bray-French K, Hartman K, Steiner G, Marban-Doran C, Bessa J, Campbell N, Martin-Facklam M, Stubenrauch K-G, Solier C, Singer T, et al. Managing the impact of immunogenicity in an era of immunotherapy: from bench to bedside. *J Pharm Sci*. 2021;110(7):2575–2584. doi:10.1016/j.xphs.2021.03.027.
- Jain T, Sun T, Durand S, Hall A, Houston NR, Nett JH, Sharkey B, Bobrowicz B, Caffry I, Yu Y, et al. Biophysical properties of the clinical-stage antibody landscape. *Proceedings of the National Academy of Sciences (PNAS)*; Vol. 114. 2017. p. 944–949.
- Jain T, Boland T, Vásquez M. Identifying developability risks for clinical progression of antibodies using high-throughput in vitro and in silico approaches. *mAbs*. 2023;15(1):2200540. doi:10.1080/19420862.2023.2200540.
- Lauer TM, Agrawal NJ, Chennamsetty N, Egodage K, Helk B, Trout BL. Developability index: a rapid in silico tool for the screening of antibody aggregation propensity. *J Pharm Sci*. 2012;101(1):102–115. doi:10.1002/jps.22758.
- Seeliger D, Tosatto SCE. Development of scoring functions for antibody sequence assessment and optimization. *PLOS ONE*. 2013;8(10):e76909. doi:10.1371/journal.pone.0076909.
- Hebditch M, Warwicker J. Charge and hydrophobicity are key features in sequence-trained machine learning models for predicting the biophysical properties of clinical-stage antibodies. *PeerJ*. 2019;7(1):e8199. doi:10.7717/peerj.8199.
- Raybould MIJ, Marks C, Krawczyk K, Taddese B, Nowak J, Lewis AP, Bujotzek A, Shi J, Deane CM. Five computational developability guidelines for therapeutic antibody profiling. *Proceedings of the National Academy of Sciences (PNAS)*. 2019;112:4025–4030.
- Negron C, Fang J, McPherson MJ, Stine WB, McCluskey AJ. Separating clinical antibodies from repertoire antibodies, a path

- to in silico developability assessment. *mAbs*. 2022;14(1):2080628. doi:10.1080/19420862.2022.2080628.
33. Kim J, McFee M, Fang Q, Abdin O, Kim PM. Computational and artificial intelligence-based methods for antibody development. *Trends Pharmacological Sci*. 2023;44(3):175–189. doi:10.1016/j.tips.2022.12.005.
 34. Bai G, Sun C, Guo Z, Wang Y, Zeng X, Su Y, Zhao Q, Ma B. Accelerating antibody discovery and design with artificial intelligence: recent advances and prospects. *Semin Cancer Biol*. 2023;95:13–24. doi:10.1016/j.semcancer.2023.06.005.
 35. Kidera A, Konishi Y, Oka M, Ooi T, Scheraga HA. Statistical analysis of the physical properties of the 20 naturally occurring amino acids. *J Protein Chem*. 1985;4(1):23–55. doi:10.1007/BF01025492.
 36. Atchley WR, Zhao J, Fernandes AD, Drüke T. Solving the protein sequence metric problem. *Proceedings of the National Academy of Sciences (PNAS)*; Vol. 102. 2005. p. 6395–6400.
 37. Sequeira AM, Lousa D, Rocha M. ProPythia: a python package for protein classification based on machine and deep learning. *Neurocomputing*. 2022;484:172–182. doi:10.1016/j.neucom.2021.07.102.
 38. Bailly M, Mieczkowski C, Juan V, Metwally E, Tomazela D, Baker J, Uchida M, Kofman E, Raoufi F, Motlagh S, et al. Predicting antibody developability profiles through early stage discovery screening. *mAbs*. 2020;12(1):1743053. doi:10.1080/19420862.2020.1743053.
 39. Fernández-Quintero ML, Ljungars A, Waibl F, Greiff V, Terje Andersen J, Gjolberg TT, Jenkins TP, Gunnar Voldborg B, Grav L-M, Kumar S, et al. Assessing developability early in the discovery process for novel biologics. *mAbs*. 2023;15(1):2171248. doi:10.1080/19420862.2023.2171248.
 40. Seeliger D, Fenn TD, Karow-Zwick AR. Developability predictions for antibody engineering and risk mitigation. *Am Pharm Rev*. 2016;19(4):188776.
 41. Jarasch A, Koll H, Regula JT, Bader M, Papadimitriou A, Kettenberger H. Developability assessment during the selection of novel therapeutic antibodies. *J Pharm Sci*. 2015;104(6):1885–1898. doi:10.1002/jps.24430.
 42. Narayanan H, Dingfelder F, Butté A, Lorenzen N, Sokolov M, Arosio P. Machine learning for biologics: opportunities for protein engineering, developability, and formulation. *Trends Pharmacological Sci*. 2021;42(3):151–165. doi:10.1016/j.tips.2020.12.004.
 43. Bepler T, Berger B. Learning the protein language: evolution, structure, and function. *Cell Syst*. 2021;12(6):654–669.e3. doi:10.1016/j.cels.2021.05.017.
 44. Lin Z, Akin H, Rao R, Hie B, Zhu Z, Lu W, Smetanin N, Verkuil R, Kabeli O, Shmueli Y, et al. Evolutionary-scale prediction of atomic level protein structure with a language model. *bioRxiv*. 2022;2022:500902.
 45. Oobatake M, Ooi T. Hydration and heat stability effects on protein unfolding. *Prog Biophys Mol Biol*. 1993;59(3):237–284. doi:10.1016/0079-6107(93)90002-2.
 46. Kozłowski LP. IPC – isoelectric point calculator. *Biol Direct*. 2016;11(1):55. doi:10.1186/s13062-016-0159-9.
 47. Wu TT, Johnson G, Kabat EA. Length distribution of CDRH3 in antibodies. *Proteins Struct Function Bioinf*. 1993;16(1):1–7. doi:10.1002/prot.340160102.
 48. Olsen TH, Moal IH, Deane CM, Lengauer T. AblLang: an antibody language model for completing antibody sequences. *Bioinf Adv*. 2022;2(1):vbac046. doi:10.1093/bioadv/vbac046.
 49. Prihoda D, Maamary J, Waight A, Juan V, Fayadat-Dilman L, Svozil D, Bitton DA. BioPhi: a platform for antibody design, humanization, and humanness evaluation based on natural antibody repertoires and deep learning. *mAbs*. 2022;14(1):2020203. doi:10.1080/19420862.2021.2020203.
 50. Ruffolo JA, Chu L-S, Pooja Mahajan S, Gray JJ. Fast, accurate antibody structure prediction from deep learning on massive set of natural antibodies. *Nat Commun*. 2023;14(1):2389. doi:10.1038/s41467-023-38063-x.
 51. Linting M, Meulman JJ, Patrick JFG, van der Kooij AJ. Nonlinear principal components analysis: introduction and application. *Psychol Met*. 2007;12(3):336–358. doi:10.1037/1082-989X.12.3.336.
 52. Cieslak MC, Castelfranco AM, Roncalli V, Lenz PH, Hartline DK. T-distributed stochastic neighbor embedding (t-sne): a tool for eco-physiological transcriptomic analysis. *Mar Genomics*. 2020;51:100723. doi:10.1016/j.margen.2019.100723.
 53. Yang Y, Sun H, Zhang Y, Zhang T, Gong J, Wei Y, Duan Y-G, Shu M, Yang Y, Wu D, et al. Dimensionality reduction by UMAP reinforces sample heterogeneity analysis in bulk transcriptomic data. *Cell Rep*. 2021;36(4):109442. doi:10.1016/j.celrep.2021.109442.
 54. Guimaraes Koch SS, Thorpe R, Kawasaki N, Lefranc M-P, Malan S, Martin ACR, Mignot G, Plückthun A, Rizzi M, Shubat S, et al. International nonproprietary names for monoclonal antibodies: an evolving nomenclature system. *mAbs*. 2022 12. 14(1):2075078. doi:10.1080/19420862.2022.2075078.
 55. Marks C, Hummer AM, Chin M, Deane CM, Martelli PL. Humanization of antibodies using a machine learning approach on large-scale repertoire data. *Bioinformatics*. 2021;37(22):4041–4047. doi:10.1093/bioinformatics/btab434.
 56. Thullier P, Huish O, Pelat T, Martin ACR. The humanness of macaque antibody sequences. *J Mol Biol*. 2010;396(5):1439–1450. doi:10.1016/j.jmb.2009.12.041.
 57. Chicco D, Jurman G. The advantages of the Matthews correlation coefficient (MCC) over F1 score and accuracy in binary classification evaluation. *BMC Genomics*. 2020;21(1):6. doi:10.1186/s12864-019-6413-7.
 58. Abhinandan KR, Martin ACR. Analysis and improvements to Kabat and structurally correct numbering of antibody variable domains. *Mol Immunol*. 2008;45(14):3832–3839. doi:10.1016/j.molimm.2008.05.022.
 59. Stewart A, C-F Ng J, Wallis G, Tsioligka V, Fraternali F, Dunn-Walters DK. Single-cell transcriptomic analyses define distinct peripheral B cell subsets and discrete development pathways. *Front Immunol*. 2021;12:602539. doi:10.3389/fimmu.2021.602539.
 60. Teixeira AAR, Erasmus MF, D'Angelo S, Naranjo L, Ferrara F, Leal-Lopes C, Durrant O, Galmiche C, Morelli A, Scott-Tucker A, et al. Drug-like antibodies with high affinity, diversity and developability directly from next-generation antibody libraries. *mAbs*. 2021;13(1):1980942. doi:10.1080/19420862.2021.1980942.
 61. Arras P, Yoo HB, Pekar L, Schröter C, Clarke T, Krah S, Klewinghaus D, Siegmund V, Evers A, Zielonka S. A library approach for the de novo high-throughput isolation of humanized VHH domains with favorable developability properties following camelid immunization. *mAbs*. 2023;15(1):2261149. doi:10.1080/19420862.2023.2261149.
 62. Brüggemann M, Osborn MJ, Ma B, Hayre J, Avis S, Lundstrom B, Buelow R. Human antibody production in transgenic animals. *Archivum immunologiae et therapiae experimentalis*. 2015;63(2):101–108. doi:10.1007/s00005-014-0322-x.
 63. Xu X, Huang Y, Pan H, Molden R, Qiu H, Daly TJ, Li N, Popoff MR. Quantitation and modeling of post-translational modifications in a therapeutic monoclonal antibody from single- and multiple-dose monkey pharmacokinetic studies using mass spectrometry. *PLOS ONE*. 2019;14(10):e0223899. doi:10.1371/journal.pone.0223899.
 64. Lin W, Wells J, Wang Z, Orengo CA, Martin ACR. Enhancing missense variant pathogenicity prediction with protein language models using VariPred. *Sci Rep*. 2024;14(1):8136. doi:10.1038/s41598-024-51489-7.
 65. Jayaram N, Bhowmick P, Martin ACR. Germline VH/VL pairing in antibodies. *Protein Eng Des Select*. 2012;25(10):523–530. doi:10.1093/protein/gzs043.
 66. Richard MG. Lessons learned from the Edrecolomab story: how a checkered past became a checkered flag for monoclonal antibodies in colorectal cancer therapy. *Oncol Res Treat*. 2005;28(6--7):311–312. doi:10.1159/000085570.

67. Stewart A, Sinclair E, Chi-Fung Ng J, Silva O'Hare J, Page A, Serangeli I, Margreitter C, Orsenigo F, Longman K, Frampas C, et al. Pandemic, epidemic, endemic: B cell repertoire analysis reveals unique anti-viral responses to SARS-CoV-2, Ebola and respiratory syncytial virus. *Front Immunol.* 2022;13:807104. doi:10.3389/fimmu.2022.807104.
68. Chothia C, Lesk AM, Tramontano A, Levitt M, Smith-Gill SJ, Air G, Sheriff S, Padlan EA, Davies D, Tulip WR, et al. Conformations of immunoglobulin hypervariable regions. *Nature.* 1989;342(6252):877–883. doi:10.1038/342877a0.
69. Pedregosa F, Varoquaux G, Gramfort A, Michel V, Thirion B, Grisel O, Blondel M, Prettenhofer P, Weiss R, Dubourg V, et al. Scikit-learn: machine learning in python. *J Mach Learn Res.* 2011;12:2825–2830.
70. Yang X, Shen Q, Xu H, Shoptaw S. Functional regression analysis using an F test for longitudinal data with large numbers of repeated measures. *Stat Med.* 2007;26(7):1552–1566. doi:10.1002/sim.2609.
71. Harris CR, Millman KJ, van der Walt SJ, Gommers R, Virtanen P, Cournapeau D, Wieser E, Taylor J, Berg S, Smith NJ, et al. Array programming with numpy. *Nature.* 2020;585(7825):357–362. doi:10.1038/s41586-020-2649-2.
72. Fay MP, Proschan MA. Wilcoxon-Mann-Whitney or t-test? On assumptions for hypothesis tests and multiple interpretations of decision rules. *Stat Surv.* 2010;4:1–39. doi:10.1214/09-SS051.
73. Johnson G, Wu TT. Preferred CDRH3 lengths for antibodies with defined specificities. *Int Immunol.* 1998;10(12):1801–1805. doi:10.1093/intimm/10.12.1801.
74. Brych SR, Gokarn YR, Hultgen H, Stevenson RJ, Rajan R, Matsumura M. Characterization of antibody aggregation: role of buried, unpaired cysteines in particle formation. *J Pharm Sci.* 2010;99(2):764–781. doi:10.1002/jps.21868.
75. Mih N, Brunk E, Chen K, Catoi E, Sastry A, Kavvas E, Monk JM, Zhang Z, Palsson BO, Valencia A. Ssbio: a python framework for structural systems biology. *Bioinformatics.* 2018;34(12):2155–2157. doi:10.1093/bioinformatics/bty077.
76. Rice P, Longden I, Bleasby A. EMBOSS: the European molecular biology open software suite. *Trends Genet.* 2000;16(6):276–277. doi:10.1016/S0168-9525(00)02024-2.