

<https://doi.org/10.1038/s42003-024-06561-3>

# Biophysical cartography of the native and human-engineered antibody landscapes quantifies the plasticity of antibody developability

Check for updates

Habib Bashour <sup>1,2,10</sup> ✉, Eva Smorodina <sup>1,10</sup>, Matteo Pariset <sup>3,10</sup>, Jahn Zhong <sup>1,4,10</sup>, Rahmad Akbar <sup>1</sup>, Maria Chernigovskaya<sup>1</sup>, Khang Lê Quý<sup>1</sup>, Igor Snapkow<sup>5</sup>, Puneet Rawat<sup>1</sup>, Konrad Krawczyk <sup>6</sup>, Geir Kjetil Sandve <sup>7</sup>, Jose Gutierrez-Marcos <sup>2</sup>, Daniel Nakhaee-Zadeh Gutierrez<sup>3</sup>, Jan Terje Andersen <sup>1,8,9</sup> & Victor Greiff <sup>1</sup> ✉

Designing effective monoclonal antibody (mAb) therapeutics faces a multi-parameter optimization challenge known as “developability”, which reflects an antibody’s ability to progress through development stages based on its physicochemical properties. While natural antibodies may provide valuable guidance for mAb selection, we lack a comprehensive understanding of natural developability parameter (DP) plasticity (redundancy, predictability, sensitivity) and how the DP landscapes of human-engineered and natural antibodies relate to one another. These gaps hinder fundamental developability profile cartography. To chart natural and engineered DP landscapes, we computed 40 sequence- and 46 structure-based DPs of over two million native and human-engineered single-chain antibody sequences. We find lower redundancy among structure-based compared to sequence-based DPs. Sequence DP sensitivity to single amino acid substitutions varied by antibody region and DP, and structure DP values varied across the conformational ensemble of antibody structures. We show that sequence DPs are more predictable than structure-based ones across different machine-learning tasks and embeddings, indicating a constrained sequence-based design space. Human-engineered antibodies localize within the developability and sequence landscapes of natural antibodies, suggesting that human-engineered antibodies explore mere subspaces of the natural one. Our work quantifies the plasticity of antibody developability, providing a fundamental resource for multi-parameter therapeutic mAb design.

Monoclonal antibodies (mAbs) are widely used therapeutics against cancer, autoimmune, and infectious diseases<sup>1–5</sup>. The global mAb market is forecasted to grow to >\$ 300 billion in 2025<sup>6</sup>. Despite their commercial success, mAb discovery remains a resource- and time-consuming process resulting in a costly and lengthy clinical approval, hindering accessibility and affordability<sup>7,8</sup>. A successful mAb molecule should not only show sufficient

affinity in its target binding profile but also exhibit a desirable “developability” profile<sup>9</sup>. The term “developability” refers to a combination of intrinsic physicochemical parameters defined as developability parameters (DPs) that relate to biophysical aspects of antibodies and their formulations—including aggregation, solubility, and stability<sup>10–14</sup>. The feasibility of an antibody candidate to successfully progress from discovery to development

<sup>1</sup>Department of Immunology, University of Oslo and Oslo University Hospital, Oslo, Norway. <sup>2</sup>School of Life Sciences, University of Warwick, Coventry, UK. <sup>3</sup>Adaptiv Biosystems, Lausanne, Switzerland. <sup>4</sup>Division of Genetics, Department Biology, Friedrich-Alexander University Erlangen-Nürnberg, Erlangen, Germany. <sup>5</sup>Department of Chemical Toxicology, Norwegian Institute of Public Health, Oslo, Norway. <sup>6</sup>NaturalAntibody, Hamburg, Germany. <sup>7</sup>Department of Informatics, University of Oslo, Oslo, Norway. <sup>8</sup>Department of Pharmacology, University of Oslo and Oslo University Hospital, Oslo, Norway. <sup>9</sup>Precision Immunotherapy Alliance (PRIMA), University of Oslo, Oslo, Norway. <sup>10</sup>These authors contributed equally: Habib Bashour, Eva Smorodina, Matteo Pariset, and Jahn Zhong.

✉ e-mail: [habib.bashour@medisin.uio.no](mailto:habib.bashour@medisin.uio.no); [victor.greiff@medisin.uio.no](mailto:victor.greiff@medisin.uio.no)

is underpinned by specific DPs, which reflect its manufacturability and druggability<sup>10,15</sup>. Thus, suboptimal developability is one of the main factors of mAbs failure in preclinical and clinical development stages<sup>16–18</sup>. Therefore, the ability to predict and prospectively design developability properties, in line with clinical and manufacturing requirements, would help by reducing the time and resources invested in developing therapeutic mAbs, thus, boosting their success rate<sup>19–21</sup>.

Traditionally, developability screening is performed through a series of laborious *in vitro* assays<sup>17,22,23</sup>. Therefore, major efforts have been invested into developing real-world-relevant *in silico* tools and machine learning (ML) algorithms that can computationally quantify or predict DP values using antibody sequence and/or structure information<sup>7,24–36</sup>. Given the current high throughput of antibody structure prediction at the repertoire scale<sup>37–39</sup>, tools for computational developability determination have become available, which can be used to identify potential design shortfalls during the development and selection of lead mAb candidates<sup>40–46</sup>.

In contrast to the design and selection of pharmaceutical mAbs, the natural (or native, used interchangeably) immune system has the capacity to “design” antigen-specific antibodies with physiologically compatible and optimized biochemical properties within days to weeks<sup>47–50</sup>. Indeed, it was previously reported that antibodies obtained via, for example, humanized mice exhibit far fewer developability risks compared to mAb candidates obtained from *in vitro* display campaigns<sup>4,17,20,51</sup>. In line with these findings, clinical mAbs have been reported to exhibit high sequence identity matches (>70%) with natural antibodies for both heavy and light chains, implying that artificially developed mAbs are not entirely dissimilar from their native counterparts, thus, highlighting the relevance of mining the natural antibody repertoires for therapeutic mAb discovery<sup>52,53</sup>. These findings motivate the embedding of therapeutic mAb sequences in the developability landscape (here defined as the multidimensional distribution of DP parameter values across DPs and antibodies) of the natural (human and mouse) antibody repertoires to assess the nativeness of their developability profiles (DPLs, a DPL is a set of DP values for a given antibody sequence/structure)<sup>54</sup>. As such, a mAb candidate with a DPL falling outside the range of its variation in natural antibodies may be assumed unnatural and, therefore, more likely to exhibit undesirable *in vivo* characteristics<sup>53</sup>. Recent studies have also pointed to the futility of attempting complete separation between natural antibodies and therapeutic mAbs based solely on DP values<sup>55</sup>. Similar findings continue to emphasize the valuable knowledge that can be harnessed from interrogating the growing sequence space of naturally sourced antibody sequences to accelerate the engineering and optimization of mAb candidates<sup>2,52</sup>.

So far, the relationship between the developability landscape of the natural antibody repertoire and that of therapeutic (or, more broadly, human-engineered) antibodies remains unclear. While not all natural antibodies may be suitable as therapeutic candidates from a developability perspective<sup>7,55</sup>, we lack a large-scale overview of sequence and structure-based natural antibody developability landscapes stratified by antibody isotype and species. Furthermore, given previous low-sample size investigations, we are unaware to what extent sequence changes affect a given DP and which sequence or structure-based design restrictions may limit all-vs-all DP optimization with a given antibody sequence. Furthermore, many studies have focused on extracting developability guidelines from a limited number of successful mAbs, considering them a “gold standard” of desired developability<sup>7,10,17,28,56</sup>. In addition, most studies have focused on a small number of DPs<sup>17,57</sup>, and apply their hypotheses to limited antibody datasets comprising of a few 100 s to 1000 s antibody sequences<sup>7,10,32,56</sup> or datasets not including patent-submitted antibodies or antibodies that failed during early clinical trials<sup>7,10,58</sup>. So far, the lack of datasets with sufficient sample size has hindered an in-depth understanding of the plasticity of the antibody developability space.

Understanding the natural antibody landscape could enable the integration of both current and prospective antibody therapeutic candidates, which will improve our interpretation of the disparities in developability between human-engineered mAbs and natural antibodies (Fig. 1). To this

aim, we have built an atlas of over two million unique native antibody sequences from human and murine heavy and light chains ( $\approx 200,000$  per isotype and chain) annotated with DPs. We predicted the 3D structure of all antibodies and calculated 40 sequence-based and 46 structure-based DPs for each antibody. Using correlation and graph theory, we identified a subset of DPs that are maximally different from one another, thus delineating a non-redundant multidimensional antibody developability space. Across all antibody isotypes, we found lower interdependency among structure-based DPs in contrast to sequence-based DPs. Notably, distinct developability landscapes emerged across species (mouse, human) and antibody chains (heavy, light). In addition, we quantified DP sensitivity by analyzing the DP value distribution of mutants with single amino acid substitutions. We also found that the values of DPs measured on the conformational ensembles of antibodies evolve throughout their molecular dynamics (MD). Regarding predictability, our analysis revealed that ML is more successful in predicting the values of sequence-based DPs, indicating a less confined design landscape for structure-based DPs. Our analysis also suggested that the observed developability spaces of human-engineered antibodies are essentially subsets of the broader natural developability space (in terms of the major principal components of variation). While our study relies on computationally predicted developability, in which experimental correspondence may vary<sup>15</sup>, it serves as a practical and real-world relevant use case of integrating and charting repertoire-wide developability to guide mAb selection and development.

## Results

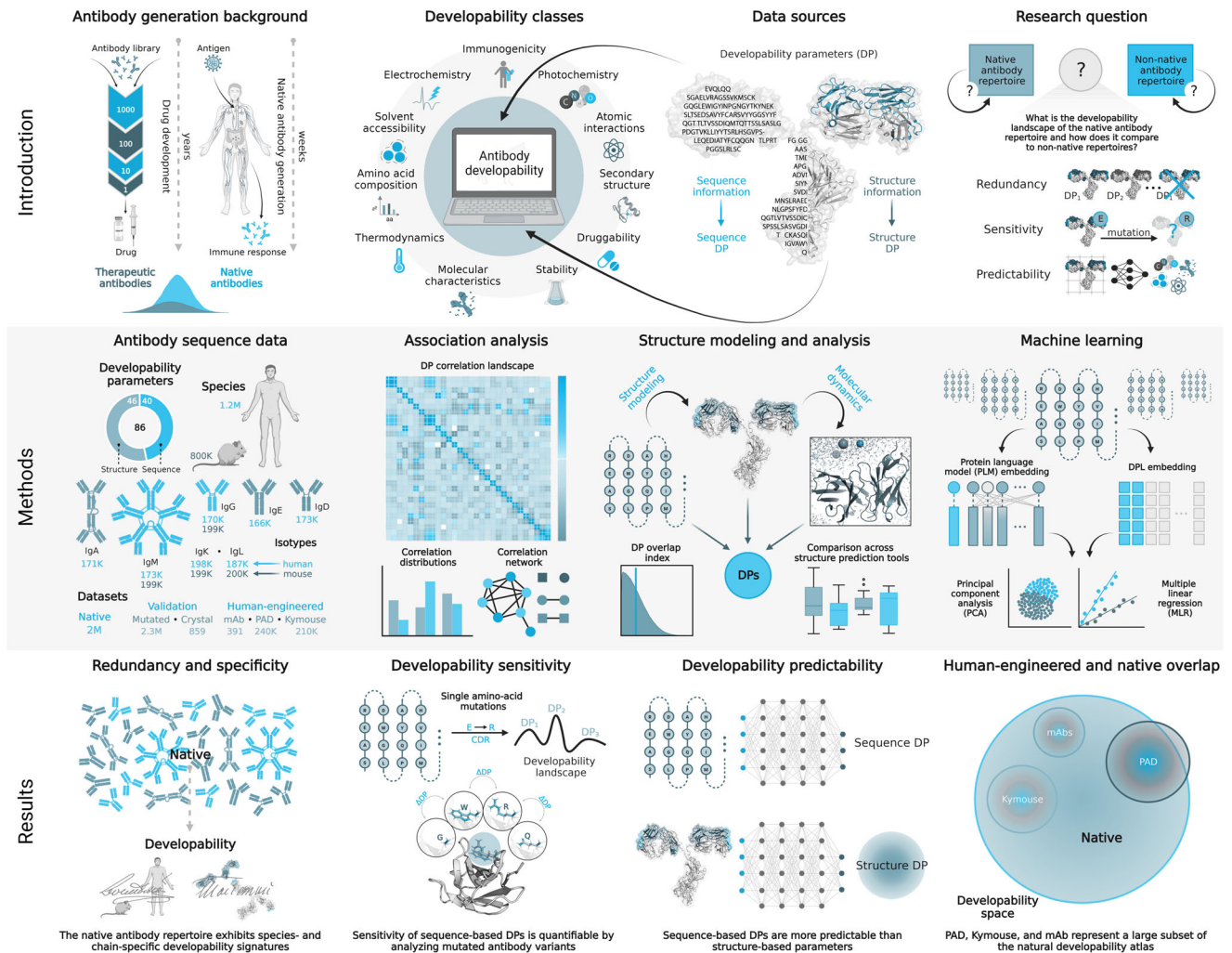
### Developability parameters

We computed 40 sequence-based and 46 structure-based developability parameters (Supplementary Data 1) for each antibody (Fv region sequence) after having predicted their 3D structures using ABodyBuilder (ABB)<sup>37</sup>. The choice of these DPs was intended to cover a comprehensive array of the physicochemical properties of antibodies based on our understanding of the developability literature and antibody structure. Some developability parameters included in the therapeutic antibody profiler (TAP, e.g. positive and negative charge heterogeneity, sequence-based hydrophobicity<sup>7</sup>) were incorporated in our study. All DPs were categorized into groups based on the main physicochemical property. For instance, the instability index and the aliphatic index were included within the “stability” group as they reflect the stability of the antibody (Supplementary Data 1). Other groups include categorical amino acid composition, electro- and photo-chemical parameters, as well as structural interactions and conformational descriptors (detailed in Supplementary Data 1).

### Sequence-based DPs show higher association and redundancies compared to structure-based DPs

The diversity and redundancy of the natural DP landscape of human and murine antibody repertoires has not been investigated. To address this knowledge gap, we assembled a dataset of  $\sim 2$  M non-paired  $V_H$  and  $V_L$  native antibody sequences ( $\sim 170K$ – $200K$  sequences for each isotype, human; IgD, IgM, IgG, IgA, IgK, and IgL, murine; IgM, IgG, IgK—Supplementary Fig. 1A).

Prior to the DP analysis, we controlled for the quality of the computational data generated along two axes Briefly, (1) given that the majority of the work was performed on unpaired chain data due to a lack of isotype-stratified paired-chain data, we verified that structure-based DPs measured on a subset of paired-chain antibodies strongly correlate with the corresponding DP values measured on each of the unpaired (heavy and light) single chains separately (median Pearson correlation 0.84–0.9, Supplementary Fig. 2, Supplementary Note 1). (2) We also analyzed, using rigid models and molecular dynamics (MD), the dependence of structure-based DPs on computational antibody structure prediction methods (such as IgFold<sup>39</sup> and ABodyBuilder<sup>38</sup>) (Supplementary Fig. 3, Supplementary Fig. 4, Supplementary Fig. 5) We found that the predominant structure prediction method used in this study (ABodyBuilder<sup>37</sup>) faithfully replicated conformations within the antibody structure conformational ensemble of a



**Fig. 1 | Redundancy, sensitivity, and predictability of antibody developability parameters in native and human-engineered antibodies.**

**Introduction:** The development of therapeutic mAbs takes years, and DPs dictate the selection and design of candidates for (pre-)clinical testing. Here, we analyzed the plasticity of the developability landscapes of natural antibodies in terms of DP redundancy (extent of DP inter-correlation), sensitivity (extent of DP change as a function of antibody sequence change), and predictability (predictability of a given DP based on one or several DPs). **Methods:** To analyze the constraints on natural antibody developability and to relate these to current human-engineered antibody datasets, we assembled a dataset of over 2 M native antibody sequences (heavy and light chain isotypes, human and murine) and computed 40 sequence- and 46 structure-based DPs. To reduce redundancy, we determined the minimum-weight dominating sets (MWDS) of DP correlation networks. To quantify sensitivity, we analyzed single-amino-acid substituted variants followed by characterization of the impact of sequence variation on DP distribution. To compute predictability and assess the interdependence of DPs, we trained multiple linear regression (MLR) using developability profile (DPL) and

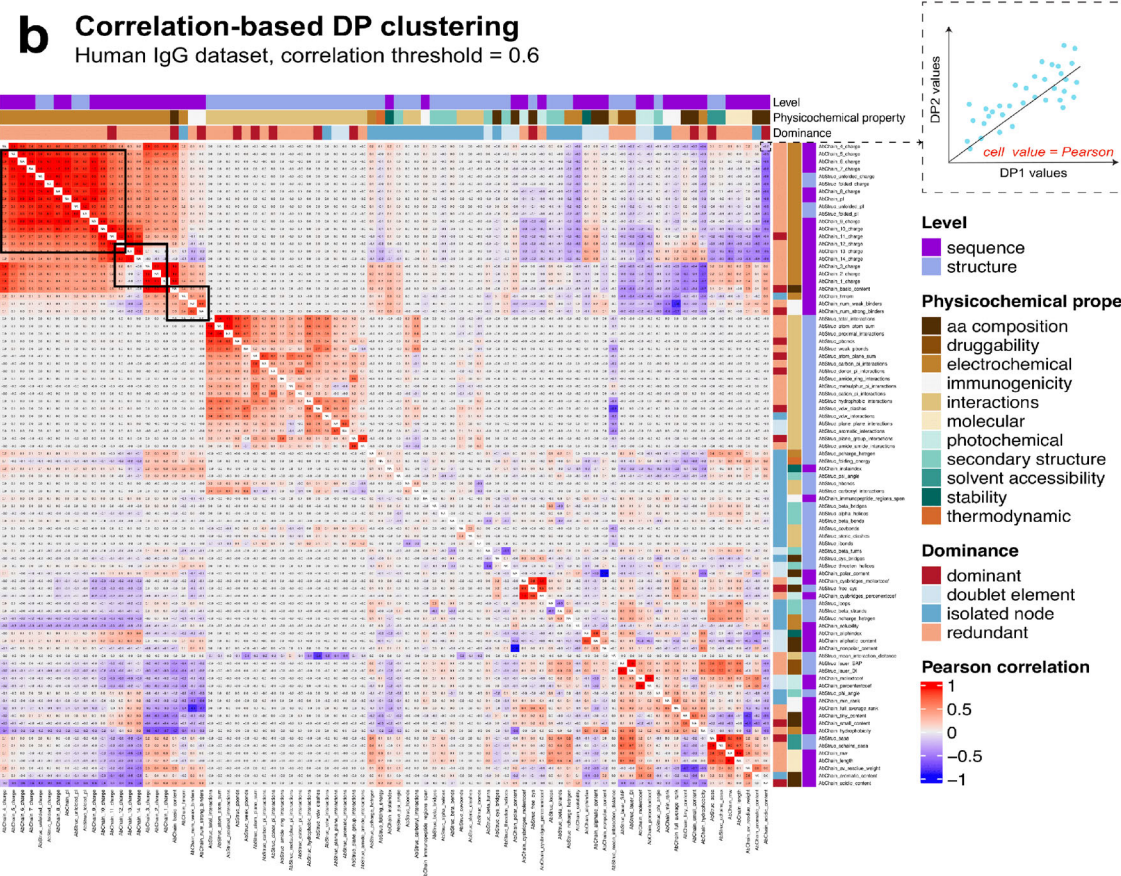
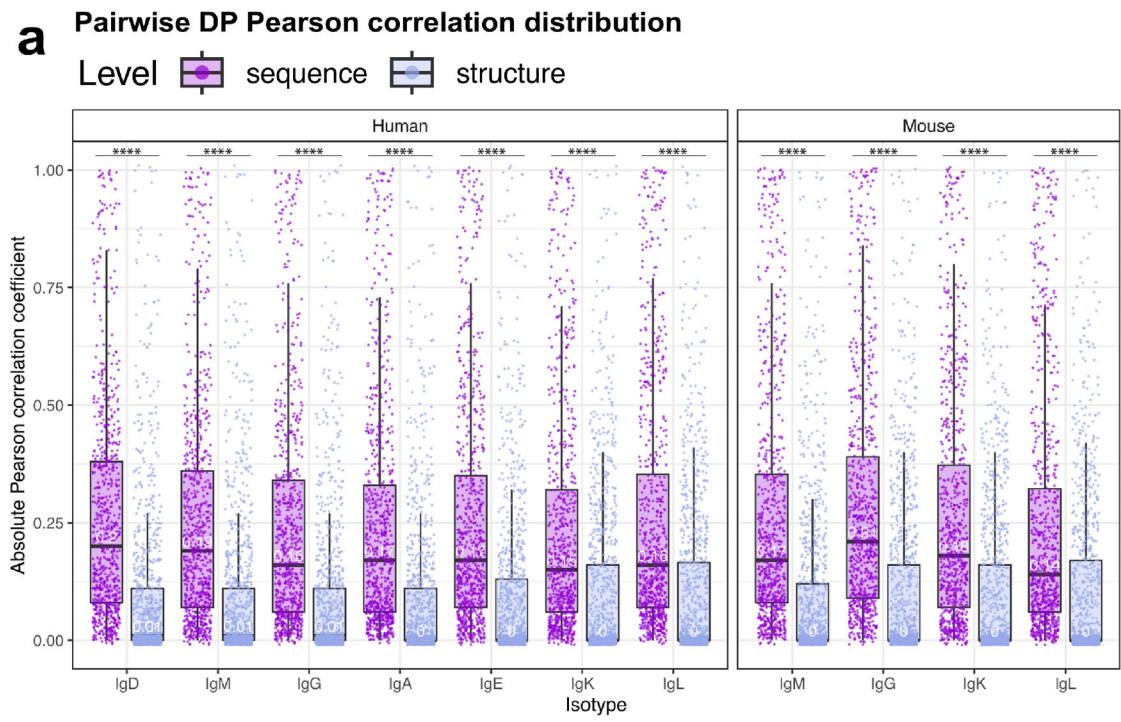
protein language model (PLM) embeddings. These embeddings were used to relate native antibodies to human-engineered ones via principal component analysis (PCA). Moreover, we performed classical molecular dynamics simulations to analyze the distributions of antibody DP values and define how the rigid models fit into these distributions. **Results:** Our results address all three research areas (redundancy, sensitivity, and predictability). **Redundancy:** We found a lower degree of inter-dependence among structure DPs than the sequence-based ones for all isotypes of the native dataset, and higher pairwise antibody sequence similarity was not always associated with higher pairwise antibody developability similarity. Native antibody datasets contained species- and chain-specific developability signatures. **Sensitivity:** We propose methods to quantify the sensitivity of antibody DPs to minimal sequence changes. **Predictability:** We found that structure-based DPs are less predictable than sequence-based DPs using protein language model (PLM) and multiple linear regression (MLR) embeddings. The comparison between native and human-engineered datasets revealed that human-engineered (therapeutic, patented, and Kymouse) datasets were localized within the native developability landscape.

given antibody sequence as determined by MD (Supplementary Fig. 4, Supplementary Fig. 5, Supplementary Note 2), thus validating our strategy to compare antibody developability landscapes.

First, we inquired about the degree of correlation among different DPs within the native antibody dataset. To address this, we examined the pairwise correlation among the values of DPs by isotype and species for the full native dataset (Fig. 2a). Overall, we found that sequence-based parameters were significantly more correlated with one another compared to structure DPs across all isotypes of the native dataset, suggesting greater general association among sequence DPs (median absolute Pearson correlation coefficient 0.14–0.21 for sequence-based DPs, 0 for structure-based DPs, Fig. 2a). Although the degree of correlation for both sequence and structure

DPs was relatively low (median absolute Pearson correlation  $\leq 0.21$ ), sequence-based parameters formed larger correlation clusters with high correlation values ( $>0.6$ ) in the native human IgG dataset (Fig. 2b; boxed in black). Similar clustering patterns with high correlation values were found across non-IgG isotypes for the sequence-based DPs (Supplementary Fig. 6, Supplementary Fig. 7). As a control measure, we repeated this analysis on permuted DPs and reported a drastic loss of correlation and a significant change in the distribution of Pearson correlation coefficient values (Supplementary Fig. 8).

Subsequently, we asked which DPs are redundant, as quantified by the previous analysis. To answer this question, we utilized the pairwise Pearson correlation coefficient values from Fig. 2a, b to construct undirected



weighted network graphs (Supplementary Fig. 9A, B). Additionally, we employed the hybrid artificial bee colony—estimation of distribution hybrid algorithm (ABC-EDA<sup>59</sup>) to identify the minimum weighted dominating set (MWDS) among DPs (see Methods). In this context, the MWDS comprises the most uncorrelated DPs that sufficiently reflect the overall developability for a set of antibodies at a given Pearson correlation threshold<sup>59</sup>. When we

conducted this analysis on the native IgG repertoire at an absolute correlation threshold of 0.6, we found that the pairwise relationships among DPs can form subnetworks, doublets, and isolated nodes (Fig. 2b, Supplementary Fig. 9B). At this threshold, subnetworks were largely formed by DPs that reflect similar physicochemical properties regardless of their level (sequence or structure—Supplementary Data 1). For instance, sequence- and

**Fig. 2 | Sequence-based developability parameters show higher redundancies compared to structure-based parameters.** **a** Absolute pairwise Pearson correlation of sequence and structure developability parameters within the native antibody dataset. Numerical values on the figure represent the median of Pearson correlation for the corresponding subset. Differences were assessed using pairwise Mann-Whitney test with  $p$ -value adjustment (Benjamini-Hochberg method). \*\*\*\* $p < 0.0001$  (Human; IgD:  $2.09e^{-112}$ , IgM:  $6.54e^{-104}$ , IgG:  $5.21e^{-100}$ , IgA:  $8.94e^{-101}$ , IgE:  $1.04e^{-94}$ , IgK:  $4.61e^{-74}$ , IgL:  $7.46e^{-80}$ , Mouse; IgM:  $1.276e^{-98}$ , IgG:  $4.76e^{-101}$ , IgK:  $6.9e^{-88}$ , IgL:  $6.68e^{-71}$ ).  $n = 785$  sequence and 1035 structure biologically independent pairwise correlation experiments for each isotype and species combination. **b** Hierarchical clustering of 40 sequence and 46 structure

developability parameters based on pairwise Pearson correlation for 170,473 IgG human antibodies (median of absolute Pearson correlation:  $0.02 \pm 0.003$  SEM). As explained in the inset (top right), each cell within the heatmap reflects the value of Pearson correlation for a pair of DPs. Developability parameters are color-annotated with their corresponding level (sequence or structure), physicochemical property (as detailed in Supplementary Data 1), and dominance status from the ABC-EDA algorithm output at Pearson correlation coefficient threshold of 0.6 (see Methods). Black boxes highlight correlation clusters that contain more than three DPs and exhibit pairwise Pearson correlation coefficient  $> 0.6$ . Supplementary Figs. 6–10A.

structure-based charge and isoelectric point (electrochemical) DPs ( $n = 20$ ) clustered with the acidic and basic amino acid composition DPs (Supplementary Fig. 9B). The largest subnetwork was predominantly occupied by sequence DPs. Meanwhile, structure-based DPs mainly formed isolated nodes ( $n = 17$ ) and smaller subnetworks (Supplementary Fig. 9B).

When we repeated this analysis on all the isotypes of the native dataset (Supplementary Fig. 1A), we found that the proportion of isolated nodes to the initial DP count was consistently higher among structure-based DPs, starting from low correlation thresholds (0.1–0.2) across all isotypes, in comparison to sequence-based DPs (Supplementary Fig. 10A). For example, only 12.5%–22.5% of sequence-based DPs (5–9 out of 40) compared to 26.1%–41.3% of structure-based DPs (12–19 out of 46) were classified as isolated nodes at a Pearson correlation threshold of 0.6 (Supplementary Fig. 10A). Meanwhile, we found that the proportion of subnetwork dominant DPs was higher among sequence DPs across all isotypes for the higher correlation thresholds (Pearson correlation  $> 0.6$ ). For example, 7.5%–12.5% of sequence-based DPs were categorized as dominant nodes at the strictest correlation threshold (0.9) in comparison to structure-based DPs (2.2%–6.5%) (Supplementary Fig. 10A). Collectively, these results emphasize the lower interdependence among structure-based DPs when compared to the sequence-based counterparts.

### The native antibody dataset exhibits chain-type and species-specific developability signatures

Next, we asked to what extent the isotypes of the native datasets are similar to one another regarding DP redundancies (MWDS parameters) and associations (DP pairwise correlations). In relation to these driving questions, we also asked to what extent natural antibodies harbor chain ( $V_H, V_L$ ) and species-specific (human, mouse) DP differences. To address these questions, we first investigated the similarities in parameter redundancies among the native dataset for a given Pearson correlation value (0.6). Specifically, we explored the pairwise intersection size of the MWDS parameters on both sequence and structure levels for the human and murine antibody datasets, featuring the common isotypes (IgM and IgG) between the two species in our dataset (Fig. 3a), and all the isotypes of the human  $V_H$  antibodies (Supplementary Fig. 10B).

This analysis revealed that both heavy (IgG and IgM) and light chain human antibody datasets displayed a larger MWDS intersection size on both sequence and structure levels than the murine counterparts (Fig. 3a). For instance, human IgM and IgG datasets shared 18 sequence DPs (86% overlap) and 29 structure DPs (90% overlap) in their MWDS sets, whereas the same heavy chain isotypes of the murine dataset shared only 12 sequence DPs (75% overlap) and 23 structure DPs (77% overlap) (Fig. 3a). Moreover, the human heavy chain dataset displayed greater or comparative MWDS intersection size even when considering all five antibody isotypes (71% overlap on sequence level, 84% overlap on structure level) in comparison to the murine heavy chain dataset (IgM and IgG only) (Fig. 3a and Supplementary Fig. 10B). Similarly, the MWDS overlap for the human light chains (IgK and IgL) was greater on both levels (15 DPs—71% overlap on sequence level and 32 DPs—97% overlap on structure level) in comparison to the mouse light chain dataset (14 sequence DPs—7% overlap and 29 structure DPs—88% overlap) (Fig. 3a). Thus, our findings suggest greater consistency among the isotypes in the human antibody dataset when it comes to DP redundancies (MWDS overlap), as opposed to the mouse antibody dataset.

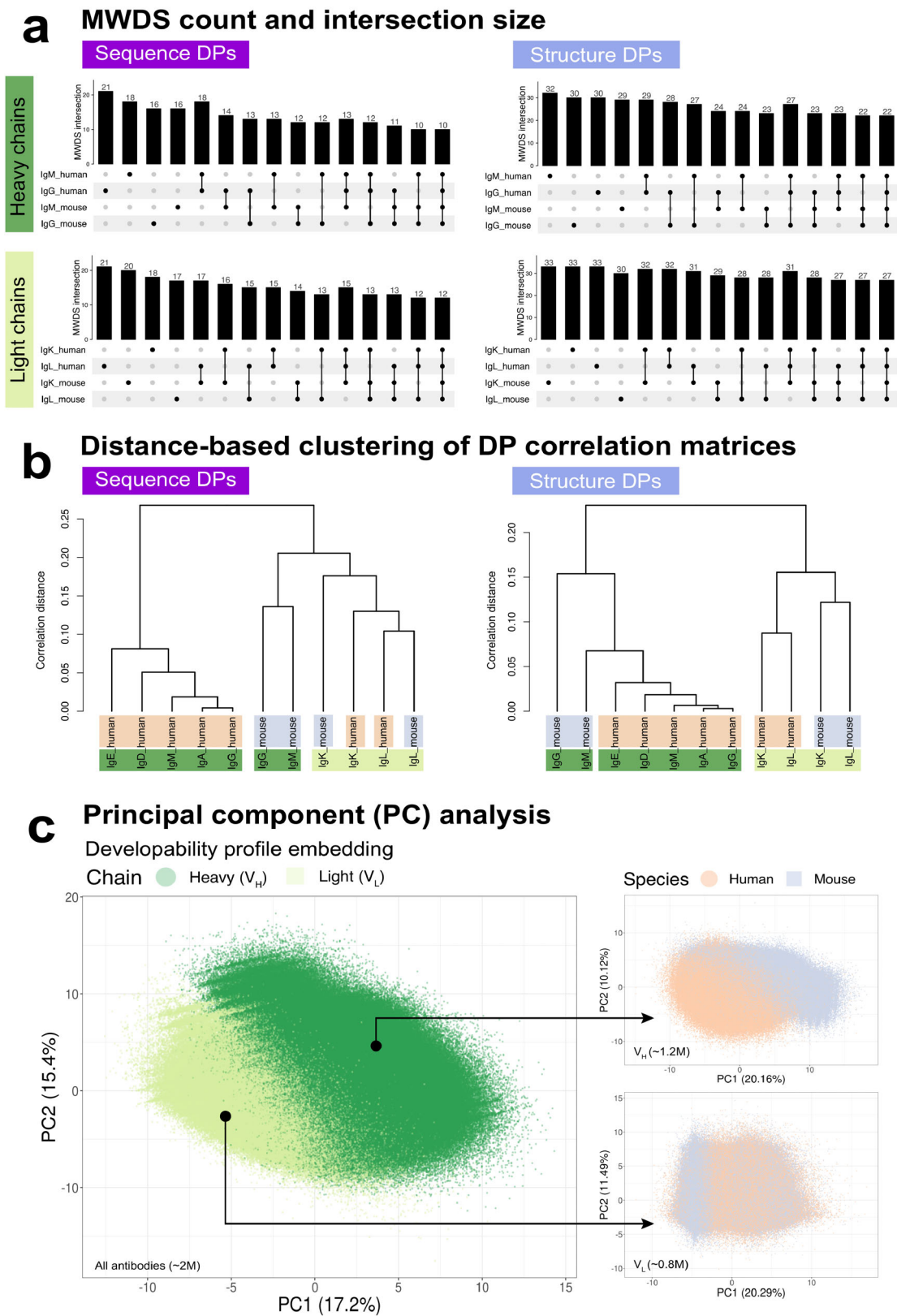
Second, we sought to investigate the similarities in DP associations among the isotypes of the native antibody dataset. To this end, we clustered the isotype-specific datasets based on the distance of their pairwise DP correlation matrices (see Methods). This analysis revealed chain-specific segregation (heavy and light) and, within a given chain, species-specific segregation (human and mouse) of antibody subsets on the structural level (Fig. 3b, right panel). Additionally, the human dataset showed a closer distance among its isotypes within the heavy and light chain clusters (0.03 for heavy chain isotypes, 0.08 for light chain isotypes) than the mouse dataset (0.15 and 0.12 for heavy and light chain clusters, respectively). An equivalent sequence-based analysis (Fig. 3b, left panel) drew a similar conclusion regarding the uniqueness of chain-type developability. However, interspecies isotype-specific clustering occurred among the light chain subsets (Fig. 3b). Similarly to the structure-based DP association clustering, the human heavy chain isotypes showed a smaller distance among themselves (0.08) in comparison to the murine IgM and IgG subsets (0.14) (Fig. 3b). These findings suggest native (human and murine) datasets harbor chain-type-specific developability signatures. Species-specific developability differences were less pronounced, especially for the light-chain antibody subsets.

The aforementioned clustering of antibody datasets based on DP associations led us to investigate whether the antibody species and chain type are key sources of variance in DP values. To this end, we performed a dimensionality reduction analysis on the developability profiles of the native antibodies using a principal component (PC) analysis (PCA) (Fig. 3c). Examining the 2D PCA projections of developability profiles of all native antibodies ( $\sim 2$  M) further emphasized differences in antibody developability by antibody chain type (Fig. 3c). The axes of maximal variance (PC1 and PC2) separated antibody sequences by  $V_H$  and  $V_L$  chain (absolute differences of medians = 5.6 and 1.8, respectively—Supplementary Fig. 10C). A similar projection of each chain type subset (heavy;  $\sim 1.2$  M, light;  $\sim 0.8$  M) allowed for (partial) species-based distinction of antibodies (Fig. 3c). The influence of the antibody species on developability was more prominent among the heavy chain antibodies (absolute difference of PC1 medians = 4.1) in comparison to the light chain ones (absolute difference of PC1 medians = 3.2) (Supplementary Fig. 10C).

In summary, the isotypes of the human native dataset exhibit greater pairwise relatedness regarding their DP associations and redundancies compared to the murine dataset. Moreover, the antibody chain type and species of origin considerably influence its overall developability.

### The sensitivity of sequence-based developability parameters is quantifiable by single-amino acid substitution analysis

Future antibody design will be performed in a multi-objective manner<sup>3,60</sup>, which means that the design approach optimizes many parameters at once. In certain cases, introducing minor changes might be sufficient to improve the value of a certain developability parameter. However, improving one parameter may compromise another<sup>3,61</sup>. Therefore, it is interesting to understand to what extent minor sequence changes impact DP values. To address this question, we performed a single-substitution sensitivity analysis of DPs by quantifying the changes in DP value induced by every possible amino acid alteration in the antibody sequence at a time (see Methods, Fig. 4a). Since there are no established methods and metrics to quantify the sensitivity of antibody developability parameters<sup>62</sup>, we employed two proxy



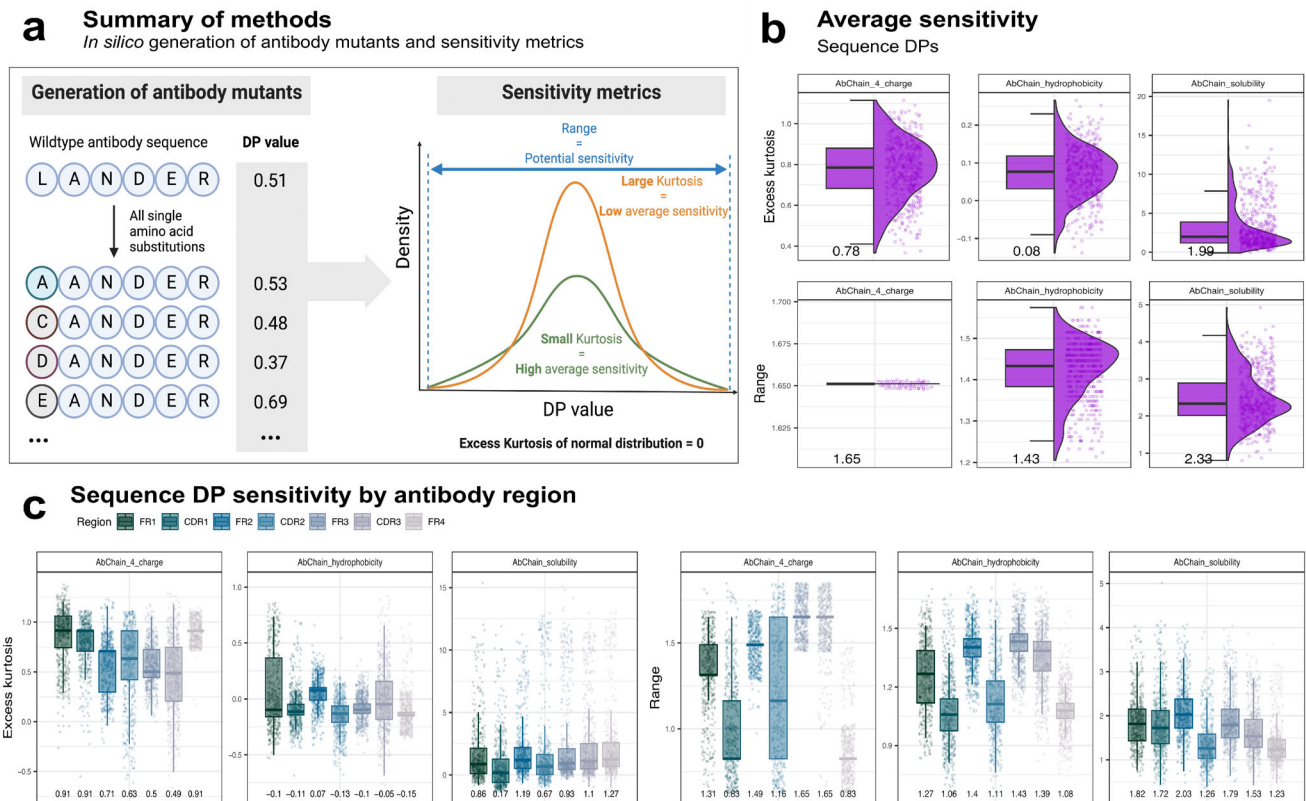
measures to estimate DP sensitivity. First, we define the *average sensitivity* by the excess kurtosis, and secondly, the *potential sensitivity* as the range of a DP distribution of an antibody and all its possible single amino acid substituted variants (see Methods). Given the inability of current antibody structure tools, both template-based or de novo deep learning-based, to

resolve differences between structures of antibodies with single amino acid variations (Supplementary Fig. 11, Supplementary Note 3), we focused our analysis on sequence-based DPs only.

We linked the dispersion of a given DP value distribution (from mutated variants and their corresponding sampled wildtype) to its

**Fig. 3 | The native (human and murine) antibody datasets exhibit chain-specific and species-specific developability signatures.** **a** MWDS intersection size for the human and mouse native datasets. Numerical values on the figure reflect the MWDS count (for an individual subset) and intersection size (for more than one subset). The MWDS for the respective isotypes was identified using the ABC-EDA algorithm (see Methods) at a threshold of absolute Pearson correlation of 0.6. For MWDS intersection size among all human heavy chain subsets, please refer to Supplementary Fig. 10B. **b** Distance-based hierarchical clustering of isotype-specific pairwise DP correlation matrices (sequence and structure levels). The height of the dendrograms

(shown to the left of the dendrograms) represents the correlation distance among the dendrogram tips. **c** Repertoire-wide principal component analysis (PCA) of the native antibody developability profiles. We performed this analysis for the complete native dataset (left pane; ~2 M sequences) and for the chain-specific datasets (right panels; ~1.2 M sequences in the top panel, ~0.8 M sequences in the bottom panel). The dimensionality of complex developability profiles was reduced to 2D PCA projections. The full value distribution of the corresponding PCs associated with each projection is shown in Supplementary Fig. 10C. Supplementary Fig. 10B, C.



**Fig. 4 | Developability parameter sensitivity can be quantified by analyzing mutated variants of wildtype antibodies.** **a** DP values were computed for all possible single amino acid substituted mutants of 500 sampled wildtype human V<sub>H</sub> antibody sequences (100 sequences sampled per isotype; *n* = 301,777 independent mutants in total). Values of each DP were scaled and mean-centered. The sensitivity was quantified for each DP by analyzing the DP dispersion of the mutants from their corresponding

wildtype. Average sensitivity was measured by excess kurtosis (small kurtosis = high average sensitivity), while potential sensitivity was measured by the range (see Methods). **b** Average and potential sensitivity of selected sequence-based DPs. **c** Average and potential sensitivity of DPs from (B) grouped by antibody region in which the mutation occurred. In both (b) and (c), numerical values on the x-axis represent the median of the corresponding sensitivity metric. Supplementary Figs. 11 and 12.

sensitivity. We employed two metrics to quantify this dispersion (see Methods). The first metric, excess kurtosis<sup>63</sup>, was implemented as a proxy measure for the *average sensitivity*. It describes how far the “tailedness” of a given distribution deviates from that of a Gaussian distribution (i.e., excess kurtosis = 0, Fig. 4a). In this context, a positive excess kurtosis indicates a low average sensitivity and a negative excess kurtosis implies a higher sensitivity with an increased proportion of mutant DP values diverging from the wildtype. Strictly, this is only valid under the assumption of a Gaussian distribution and should be considered when estimating sensitivity by excess kurtosis.

The second metric is the range, which is the absolute difference between the smallest and largest values of a distribution after normalization. It was used as a proxy measure for the *potential sensitivity* of DPs as it reflects the potential extreme changes that can be introduced on DP values as a factor of amino acid sequence change (Fig. 4a). Since only single amino acid substitutions were analyzed, we removed all DPs describing categorical amino acid composition (Supplementary Data 1) where it is trivial to predict

the changes in DP values, as well as the length of the sequence (AbChain\_length) since it is not affected by substitutions. Because DPs denoting a sequence’s charge at a given pH were clustered in three highly correlated clusters (Fig. 2b), we chose to retain only three of them, which represent acidic, neutral and basic pH (AbChain\_4\_charge, AbChain\_7\_charge and AbChain\_12\_charge respectively). Additionally, all DPs with ‘\_percentcoef’-suffix were analyzed instead of their highly correlated counterparts with ‘\_molextcoef’-suffix (Supplementary Data 1).

We found the median excess kurtosis of most sequence-based DPs was >0 and that most DPs exhibit an *average sensitivity* close to that of a normal distribution (excess kurtosis of a normal distribution = 0, Fig. 4a, Supplementary Fig. 12A). In fact, some parameters, such as the molar extinction coefficient (of cysteine bridges) and the hydrophobic moment, were insensitive on average to substitutions as indicated by their high kurtosis (median excess kurtosis: 6.7, 49.02, and 29.49, respectively; Supplementary Fig. 12). Notably, none of the tested DPs displayed high average sensitivity (median excess kurtosis « 0, Fig. 4a, Supplementary Fig. 12), suggesting that

developability is relatively stable to the average single amino acid mutation with few outliers. Nevertheless, since DP values were normalized, small relative shifts induced by a mutation may still have a large effect in practice.

Next, we show that the median range in sequence-based DPs varied between 0.38 for the molecular weight DP (AbChain\_mw) and 3.82 for the hydrophobic moment DP (AbChain\_hmom—Supplementary Fig. 12B). Notably, some DPs (such as AbChain\_4\_charge—Fig. 4b first column second row) have a constant or close to constant range, likely due to the fact that the set of all possible single amino acid substitutions covers the entire range of possible DP values. For example, when considering the charge of an antibody, the lowest possible charge results from substituting the most positively charged amino acid with the most negatively charged amino acid and vice versa. Since all amino acids are present in close to all sampled wildtype sequences, their ranges are almost identical as well. In DPs that depend on non-linear amino acid interactions, such as solubility (AbChain\_solubility), the range was more diverse (2.33, Fig. 4b).

To investigate the impact of substitutions across antibody regions, we grouped DP values of human heavy chain mutants by the region in which the substitution occurred and calculated their sensitivity metrics separately. We observed that electrochemical DPs such as charge and hydrophobicity (AbChain\_4\_charge and AbChain\_hydrophobicity) exhibited higher potential sensitivity (range) in CDR3 and framework regions (median ranges AbChain\_4\_charge and AbChain\_hydrophobicity respectively; CDR3: 1.65, 1.39, FR1: 1.31, 1.27, FR2: 1.49, 1.4 and FR3: 1.65, 1.43) compared to CDR1 (0.83, 1.06) and CDR2 (1.16, 1.11) and FR4 (0.83, 1.06; Fig. 4c). Although we found the average sensitivity (excess kurtosis) to differ by antibody region (Fig. 4c), there was no apparent general rule that clearly separates framework regions from CDRs. Since short sequences have a higher probability of missing the most charged or polar residues, the potential sensitivities of a given region tend to be lower for shorter regions. The stark differences in range between heavy chain framework regions 1–3 and CDRs 1–2 may thus be explained by the shorter sequence length of CDRH regions (with the exception of CDRH3).

In summary, our sensitivity analysis suggests that most DPs are comparably ‘normally’ sensitive (close to 0 excess kurtosis: mutant DP distribution as ‘tailed’ as a Gaussian distribution), and some parameters are especially *insensitive* to the average substitution. Additionally, although average and potential sensitivity differ by antibody region in which a mutation occurs, the differences are not generalizable across DPs.

### Antibody sequence similarity does not imply antibody developability similarity

Given that the values of DPs were prone to change with minor sequence changes, we asked to what extent pairwise sequence similarity is related to pairwise developability profile similarity, where the developability profile (DPL) was defined as a numerical vector that carries (sequence and/or structure) DP values in a fixed order for a given antibody sequence (see Methods).

To this end, we first examined the pairwise correlation of antibody developability profiles (developability profile correlation: DPC) alongside the pairwise sequence similarity score (see Methods) for a random sample of 100 natural antibodies from the human IgM dataset that share the IGHV gene family annotation (Fig. 5a). This is to eliminate the role of the V-gene as a factor of variance in our analysis, as up to 80% of sequence similarity can be expected among antibodies that belong to the same IGHV gene family<sup>64,65</sup>. Sequence-level DPC clusters were often, but not always, accompanied by sequence similarity clusters, while structure-level DPC clusters were independent regarding sequence similarity clusters (Fig. 5a, Supplementary Fig. 13, Supplementary Fig. 14 and Supplementary Fig. 15). To quantify the association between the two metrics (DPC and sequence similarity), we computed the Pearson correlation coefficient between the pairwise DPC matrices and the pairwise sequence similarity matrices (Fig. 5b, Supplementary Fig. 16A). We repeated this analysis for 100 randomly sampled sets of 100 sequences each (within the same IGHV gene family). Samples were taken from all isotypes of the native dataset to account for variation of

associations among batches. We restricted the single set (batch) size to 100 antibodies to ensure correlation matrix regularization<sup>66</sup>. We examined the resulting Pearson correlation values alongside the average sequence similarity score (100 values for each metric for the 100 sampled sets per isotype—Fig. 5b). We found that Pearson correlation coefficients (between DPC and average sequence similarity) tended to be higher on the sequence level (0.2–0.7) than on the structure level (0.1–0.4) across all antibody isotypes (Fig. 5b). For instance, the mean Pearson correlation coefficient for the human IgD dataset was 0.5 on the sequence level and 0.2 on the structure level (Fig. 5b). This finding suggested that similar sequences exhibit higher sequence-based developability similarity compared to structure-based developability similarity. However, higher sequence similarity was not always accompanied by higher Pearson correlation values of DP profiles. For example, although the average sequence similarity of the murine IgL dataset was as high as 0.9, the mean value of Pearson correlation coefficient was only 0.5 (Fig. 5b). We reported a similar Pearson correlation average (0.5) for the human IgE dataset, even though its mean sequence similarity was less than the IgL murine dataset (0.7, Fig. 5b).

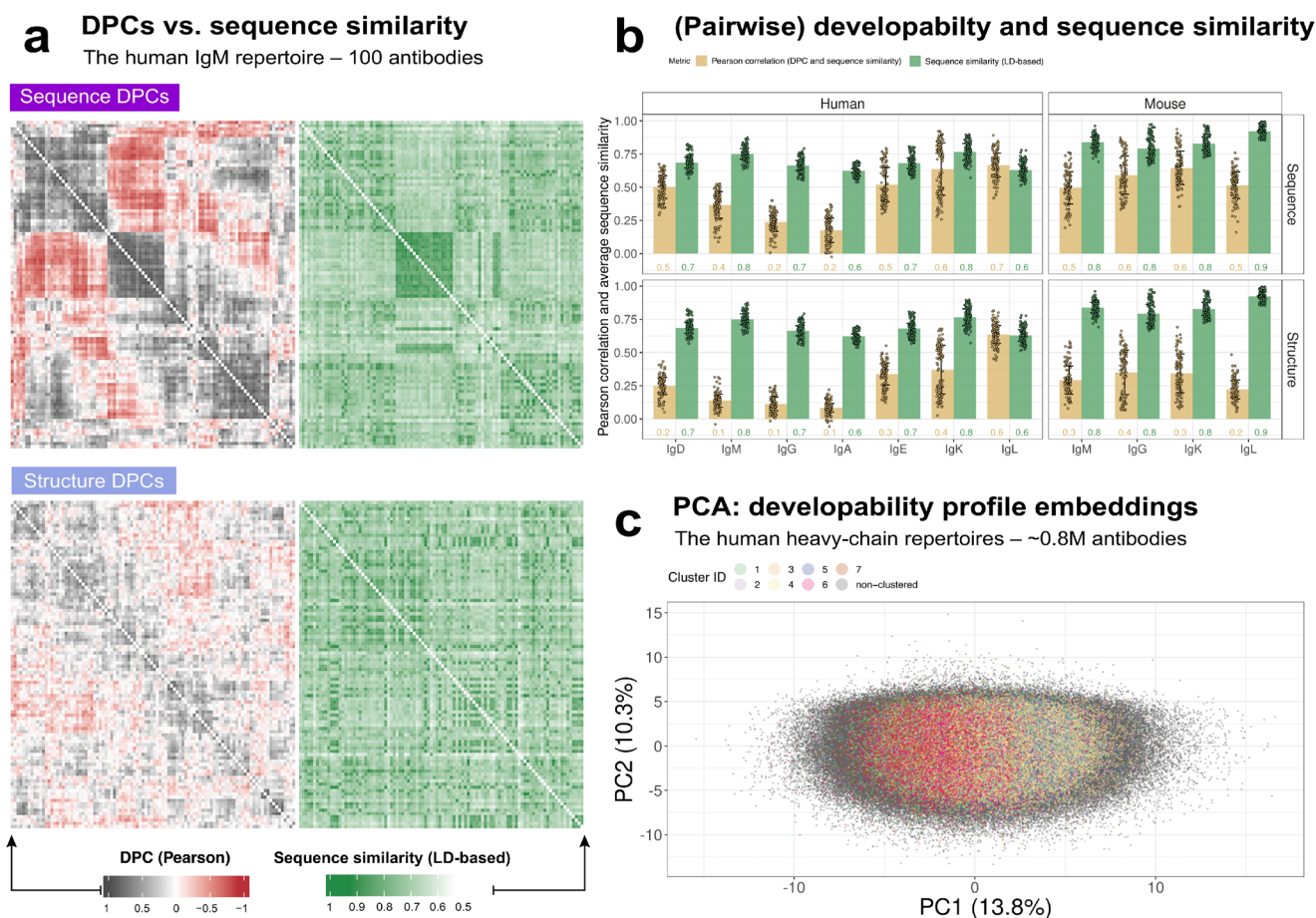
Next, we investigated the relationship between antibody developability and sequence similarity using a geometric approach to test the finding that developability profile similarity and sequence similarity are not necessarily associated (Fig. 5a, b). We leveraged the fact that antibody developability profiles are numerical vectors in the developability space ( $R^N$ ), where  $N$  represents the number of DPs that compose a single developability profile (see Methods). This space ( $R^N$ ) offers a natural notion of antibody developability diversity. However, it is challenging to inspect due to its high dimensionality. Thus, we applied a dimensionality reduction technique (principal component analysis: PCA) on the developability profile space of the human heavy-chain antibody dataset (Fig. 5c) as it is the largest data subset that belongs to a single species and chain type (~0.8 M antibodies, Supplementary Fig. 1A). Within this subset, we identified seven sequence similarity groups (1–7). Each group encapsulates at least 10 K antibodies and group members exhibit at least 75% pairwise sequence similarity (Supplementary Fig. 16B—see Methods). We found that antibody members that belong to the same sequence similarity group did not occupy a restricted space on the PCA projection plane, suggesting that antibody developability and sequence similarity are not correlated (Fig. 5c). To quantify this observation, we studied the correlation between the pairwise Euclidean distances in the developability space ( $R^N$ ) and the pairwise (normalized) Levenshtein distances for 5000 antibody sequences from the human IgM dataset that share the same IGHV gene family annotation (IGHV1) and belong to the same sequence similarity group (group 1—Supplementary Fig. 16B, C). We found that the two distance measures showed only minimal correlation (Pearson correlation coefficient = 0.18, Supplementary Fig. 16C). As Fv sequences belonging to the same IGHV gene family share high (up to 80%) sequence similarity<sup>64,65</sup>, most of which is attributed to conserved sequences in the frameworks, CDRH1 and CDRH2 regions<sup>67,68</sup>, the majority of sequence variation is credited to the CDRH3. As such, the lack of correlation between the pairwise Normalised Levenshtein distance and the pairwise Euclidean distance in the developability space for these antibodies (Supplementary Fig. 16C) indicates an independence of developability profile and CDRH3 sequence similarity among antibodies of the native dataset.

In conclusion, our analysis demonstrated that antibody developability and sequence similarity were largely independent, suggesting that improving the developability profile for a certain therapeutic mAb candidate with a desirable target binding profile may be possible by introducing small changes in its amino acid sequence.

### DP predictability implies interdependence and antibody design space restriction

Building on the prior finding that no significant association exists between antibody sequence similarity and developability similarity, we inquired whether missing values of given DPs can be predicted based on the knowledge of the values of other DPs. This is to investigate to what extent the





**Fig. 5 | Developability profile similarity is not necessarily associated with sequence similarity.** a Pairwise developability profile Pearson correlation (DPC—left panels) alongside the pairwise Levenshtein distance (LD)-based-sequence similarity score (right panels—see Methods) for a random sample of  $n = 100$  antibodies from the human IgM dataset ( $100 \times 100$  matrices) that share the same IGHV gene family (*IGHV1*) annotation (shown both for sequence and structure DPLs). Each row and each column represent a single antibody sequence. Rows and columns in the left panels were hierarchically clustered. In the right panels (sequence similarity), rows and columns were ordered in the same order as the corresponding left panel (DPC) for ease of comparison. The distribution of DPC and sequence similarity is shown in Supplementary Fig. 16A. b Pearson correlation between DPC and sequence similarity matrices for 100 sets of randomly sampled non-overlapping 100 antibody sequences (within the same IGHV gene family per batch) from all isotypes of the native dataset (total  $n = 100$  independent experiments of 100 antibodies per

experiment). Pearson correlation coefficient values (shown in beige) are presented alongside the corresponding mean sequence similarity values (shown in green) for the same 100 sets. The height of the bars and the numerical values on the figure reflect the mean of the corresponding metric (mean Pearson correlation and the mean sequence similarity). The error bars represent the standard deviation. c Principal component analysis (PCA) of the developability profiles of the native human heavy-chain dataset ( $n = \sim 0.8$  M antibodies). The developability profiles (DPLs) were utilized as embeddings for this analysis (see Methods). Antibody clusters (1–7) were created for the groups of antibodies that are at least 75% similar in sequence (as determined by USEARCH) and contain at least 10 K antibodies. Antibodies that did not satisfy the clustering conditions were labeled as “non-clustered” (727861 sequences) and sent to the back layer of the figure. For antibody counts per cluster, please refer to Supplementary Fig. 16B. Supplementary Figs. 13–16.

developability space is amenable to orthogonal DP design (Fig. 6a). In this context, high predictability of a given DP would indicate a restriction of the antibody design space, while low predictability could signify a more plastic space with a higher degree of freedom for the values of this DP. Answering the above question would also provide insights into which DPs can be better predicted (with the provided knowledge of the remaining DPs), which may accelerate antibody developability screening. Also, we evaluated the predictability of missing DP values depending on the sole knowledge of the amino acid sequences of antibodies (through their protein language model representations, Fig. 6a). This aims to investigate the feasibility of DP predictability in the absence of other DP data.

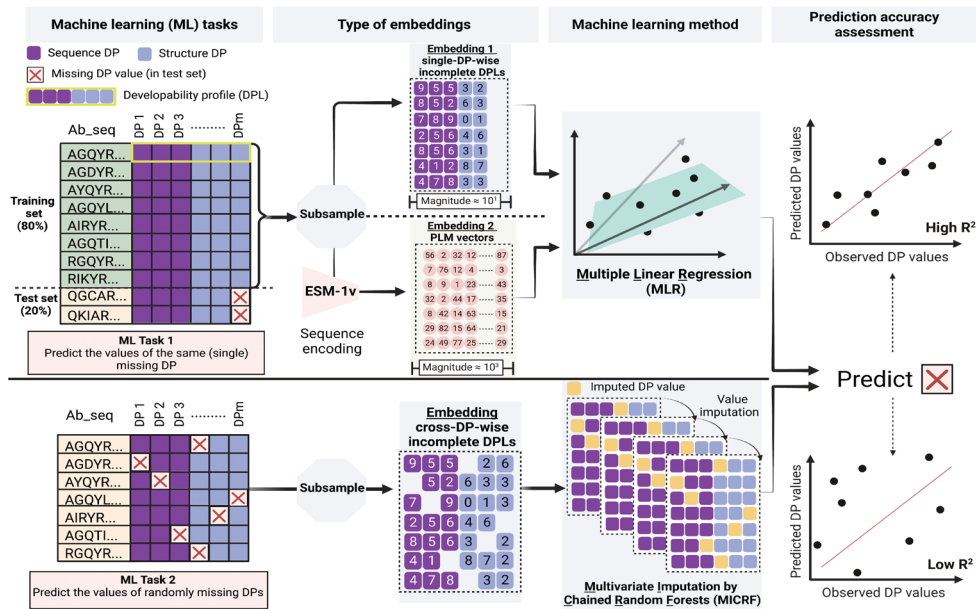
Importantly, the primary objective of this analysis is to evaluate how various factors (type of ML input, type of predicted DP, the redundancy state of the predicted DPs) affect the predictability of DPs rather than achieving absolute predictive precision given that we did not perform a comprehensive benchmarking of encoding or ML approaches<sup>69,70</sup>. Ultimately, when comprehensive benchmarking and fine-tuning of multitask ML models is performed on developability data, it might be possible to

depend on these models for DP value predictability rather than relying on several in silico tools to achieve the same task.

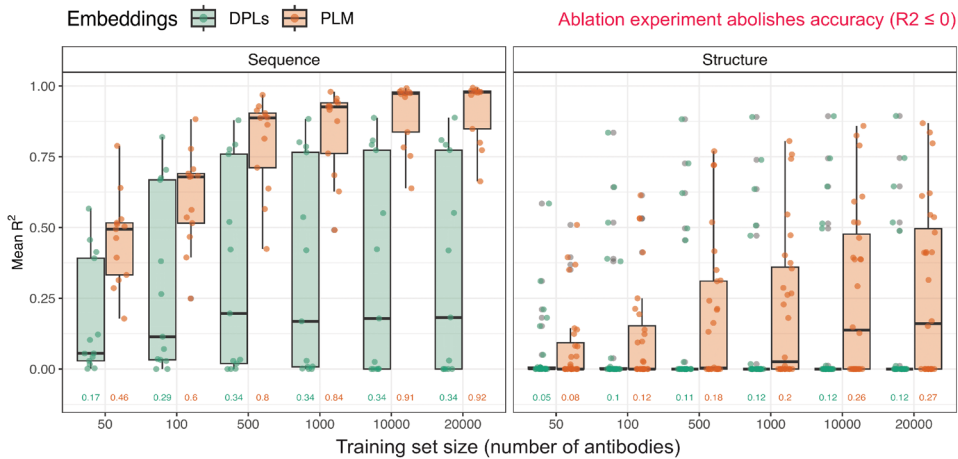
We initially used MWDS DPs to eliminate collinear dependence that may increase ML prediction accuracy in a non-interesting or trivial manner. To this end, we investigated two scenarios where the missing (deleted) DP values were either all from a single DP (ML Task 1), or randomly missing from several DPs (ML Task 2) (Fig. 6a). In practical terms, ML Task 1 replicates a use case where the values of a single DP—that might be challenging to compute or measure—are missing from all antibodies in the dataset. Meanwhile, as “spotted” (i.e., missing-at-random) data is a real-world problem in biomedical research<sup>71–73</sup>, ML Task 2 replicates the use case of obtaining a developability dataset where the values of several DPs are sporadically missing from the antibodies (Fig. 6a).

For ML Task 1, we compared the predictive accuracy of two types of (input) embeddings to predict the missing DP values via multiple linear regression (MLR) models after defining training and test (sub)sets from the native human  $V_H$  antibody dataset (Fig. 6a, see Methods). (i) The first embedding is the single-DP-wise incomplete developability profiles

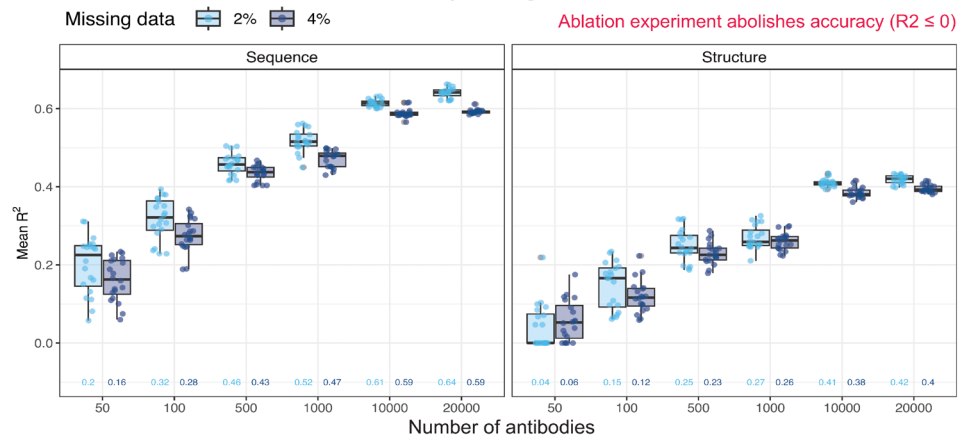
**a Graphical summary of machine learning (ML) approaches for DP value prediction**



**b ML Task 1: Predict the values of the same (single) missing DP**



**c ML Task 2: Predict the values of randomly-missing DP values**



(DPLs) which depend on the knowledge of all other DP values included in the training set (low-dimensional embedding—order of magnitude  $\approx 10^1$ —Fig. 6a). (ii) The second embedding is an amino acid sequence encoding output produced by the protein language model (PLM) ESM-1v, which generates large semantically-rich digital

representations of antibodies (high-dimensional embedding—order of magnitude  $\approx 10^3$ )<sup>74–76</sup>. Unlike DPL-based embedding, PLM-based embedding is entirely unaware of antibody DP values (Fig. 6a) and we aim to test whether these biochemical properties are implicitly contained in it.

**Fig. 6 | Sequence-based developability parameters are more predictable than structure-based parameters.** **a** Graphical representation of machine learning (ML) approaches used to assess the predictability of DPs. We investigated two scenarios where the missing (deleted) DP values were either all from one (single) DP (ML Task 1) or were randomly missing from several DPs (ML Task 2). For ML Task 1, we compared the predictive accuracy of two different embeddings; single-DP-wise incomplete developability profiles (DPLs) (embedding 1; order of magnitude  $10^1$ ) and PLM vectors (embedding 2; order of magnitude  $10^3$ ). We used these embeddings to train multiple linear regression (MLR) models (separately) to predict the missing DP values in the test set. To enable the comparison between these two embeddings, we used identical training subsamples (in regards to size and antibody identity, see Methods). For ML Task 2, we used cross-DP-wise incomplete developability profiles as input for the multivariate imputation by chained random forests (MICRF) algorithm to predict missing DP values. For both ML tasks, we estimated the prediction accuracy by computing the coefficient of determination ( $R^2$ ) using observed and predicted DP values<sup>171</sup>. **b** Comparison of the predictive accuracy of incomplete developability profiles (single-DP-wise incomplete DPLs) and PLM vectors as embeddings for MLR models to predict the values of missing DPs in the test set (ML

Task 1). The x-axis reflects the number of antibody sequences (sample size) used for the embedding. For each sample size, we repeated the prediction of missing DPs 20 times ( $n = 20$  independent experiments). The y-axis represents the mean  $R^2$  for sequence DPs (left facet) and structure DPs (right facet). Error bars represent the standard deviation of  $R^2$ . Missing DPs tested in this analysis belonged to the MWDS exclusively, as determined at a Pearson correlation coefficient threshold of 0.6, for the human IgG dataset, summing to 13 sequence DPs and 28 structure DPs (after removing a single element from each doublet and immunogenicity DPs, Supplementary Table 2). **c** Evaluating the predictability of randomly missing DP values using the MICRF algorithm where cross-DP-wise incomplete developability profiles are used as embeddings. The x-axis reflects the number of antibodies (sample size) used for the embedding. For each sample size, we repeated the prediction of missing DPs 20 times ( $n = 20$  independent experiments). The y-axis represents the mean  $R^2$  for sequence DPs (left facet) and structure DPs (right facet) when the proportion of the missing data is either 2% (light blue line) or 4% (dark blue line). Missing DPs tested in this analysis belonged to the MWDS, analogously to (b). Numbers on the x-axis in both (b) and (c) reflect the average values of mean  $R^2$ . Supplementary Fig. 17.

We found that the predictive accuracy of MLR models increased with increasing training set size for both embeddings. However, DPL-based models reached their saturation point earlier than PLM-based ones for both sequence and structure DPs (1000 antibodies for DPLs, 20000 for PLM—Fig. 6b). This is also (partly) attributable to the fact that higher dimensional inputs (PLM embeddings) correspond to additional degrees of freedom when training the model. Overall, both embeddings achieved higher prediction accuracy for sequence DPs compared to structure DPs at their saturation points (Fig. 6b). However, the disparity in prediction accuracy between the two embeddings was more pronounced for sequence DPs, with PLM-based embeddings achieving a mean prediction accuracy of 0.92 compared to 0.34 for DPL-based ones (Fig. 6b). Thus, the high predictability enabled by PLM-based embedding on sequence level highlighted its capacity to capture the biophysical properties of antibodies based on amino acid sequence<sup>77</sup>.

When conducting an analogous analysis including non-MWDS (redundant) DPs, we found notable improvement in DPL-based MLR models to predict missing DP values (mean  $R^2$  of 0.88 for sequence DPs and 0.36 for structure DPs) (Supplementary Fig. 17A). This is due to the inherent collinearity among non-MWDS DPs (Fig. 2b) that simplifies DPL-based prediction, resulting in higher prediction accuracy<sup>78</sup>. In contrast, PLM-based embeddings showed far lesser to no distinct improvements when predicting non-MWDS DPs (mean  $R^2$  of 0.96 for sequence DPs and 0.38 for structure DPs) (Supplementary Fig. 17A).

For ML Task 2, we implemented the multivariate imputation by chained random forests (MICRF) algorithm<sup>72,79</sup> to evaluate its prediction accuracy to recover randomly missing (deleted) DP values from the developability dataset (Fig. 6A—see Methods). We investigated two cases where either 2% or 4% of DP values are missing from either sequence or structure MWDS parameters (Fig. 6C), as implementing the MICRF algorithm with a greater fraction of missing data could compromise the accuracy and reliability of the imputed (predicted) DP values<sup>79</sup>. Overall, and similarly to the observations reported from ML Task 1 (Fig. 6b), structure DPs were more challenging to predict, and a larger number of data inputs (antibodies) aided the achievement of higher prediction accuracies (Fig. 6c). However, we noticed (only) subtle differences in the prediction accuracy (mean  $R^2$ ) when comparing the algorithm capacity to restore 2% or 4% of missing DP values (Fig. 6c), outlining its robustness within the advised data loss limits for its application<sup>80,81</sup>. For example, we reported mean  $R^2$  of 0.64 (sequence DPs) and 0.42 (structure DPs) with 2% data loss compared to 0.59 (sequence DPs) and 0.4 (structure DPs) with 4% data loss, at a saturation point of 20000 antibodies (Fig. 6c).

Similarly to ML Task 1, non-MWDS DPs were shown to be easier to predict (Supplementary Fig. 17B). However, the disparity in prediction accuracy between the two classes of DPs (MWDS and non-MWDS) was less pronounced in comparison to DPL-based predictions in ML Task 1. For instance, at 2% data loss, the improvement in prediction accuracy of

sequence DPs was (only) ~0.2 higher (0.85 for non-MWDS, 0.64 for MWDS—Supplementary Fig. 17B) compared to ~0.6 in DPL-based predictions in ML Task 1 (0.88 for non-MWDS, 0.34 for MWDS Supplementary Fig. 17A).

Of note, we performed ablation studies<sup>82</sup> on both ML tasks, by randomly permuting the values of DPs at the columns in the input data (feature shuffling), and confirmed that the prediction accuracy was diminished ( $R^2 \leq 0$ , Fig. 6b, c, see Methods).

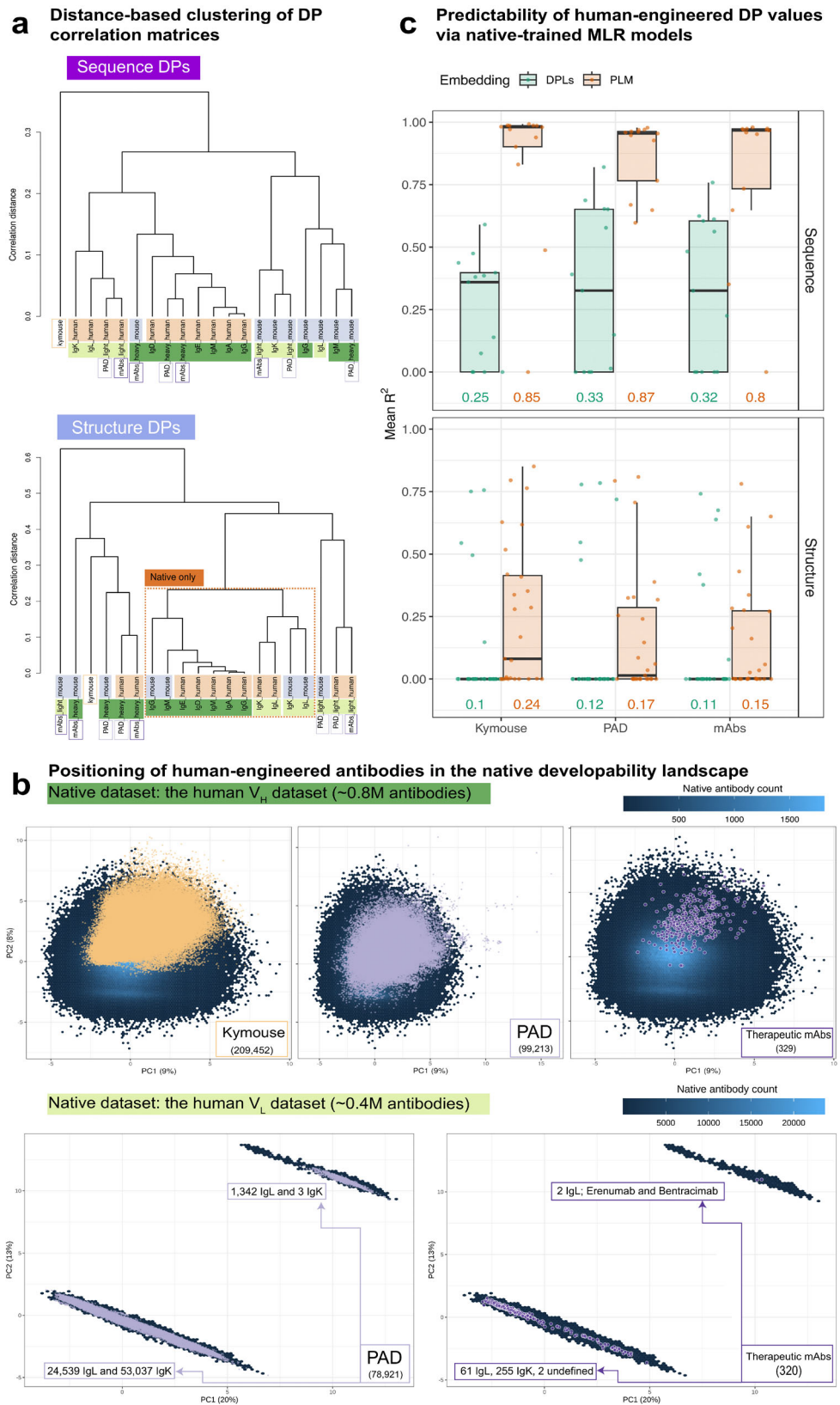
In summary, for the ML methods and embeddings applied in this analysis, we found that the structural developability space is less restricted. Additional analyses suggested that variations in structure prediction alone cannot explain this (Supplementary Notes 2–3). Additionally, this analysis highlighted the potential of PLMs to reflect the sequence-based aspects of antibody developability, if provided with sufficient training data. We want to emphasize that our main goal was not obtaining perfect prediction accuracy of DPs but rather measuring the amount of readily available biological information contained in the two representations examined, i.e. DPL- and PLM-based representations. Finally, we also addressed the presence of potential gaps in developability information in practical settings by showing how two different ML algorithms can be utilized in two different scenarios of missing data.

### Patent-submitted, humanized mouse (Kymouse) and therapeutic monoclonal antibodies represent a subset of the natural developability atlas

Previously, clinically approved therapeutic antibodies were used to build classifiers for optimal developability<sup>7,36</sup>. However, the higher abundance of native antibody data incentivizes the investigation of how human-engineered antibodies relate to their native counterparts. Thus, we investigated to what extent we can detect differences between native and human-engineered antibodies (therapeutic mAbs, Kymouse and PAD). To this end, we examined the proportional abundance of liability sequence motifs and the representation of developable germlines across the native and human-engineered datasets (Supplementary Fig. 18). Liability motifs are short amino acid sequences, which may negatively impact various aspects of antibody developability when present in their CDRs<sup>83</sup>. Developable germlines represent a group of human immunoglobulin genes ( $V_H$  and  $V_L$ ), which have been suggested to harbor favorable biophysical properties<sup>84,85</sup>.

The native dataset was comparable to the human-engineered ones regarding the proportion of antibodies that exhibit liability motifs, and no consistent increased or decreased representation (trend) of these motifs was reported (Supplementary Fig. 18A). For instance, while the native human  $V_L$  dataset exhibited a higher abundance of the asparagine deamidation motif “NS” (17.8%) compared to PAD and mAbs (15.7%, 16.6%), it exhibited a lower abundance of the proteolysis motif “DP” (0.5%, 2%, 3.8%—Supplementary Fig. 18A). Similarly, the native dataset harbored a comparable representation of the developable germline genes with the exception of *IGHV1* members (5.2% in native, 3.1% in Kymouse, 16.6% in PAD, and

**Fig. 7 | Human-engineered antibodies are contained in the developability landscape of natural antibodies.** **a** Distance-based hierarchical clustering of isotype-specific pairwise DP correlation matrices (sequence and structure levels—similar to the analysis shown in Fig. 3a). The height of the dendrograms (shown to the left of the figure) represents the correlation distance among the dendrogram tips. The dashed square in the right (structure-based) panel highlights the native-only dataset. **b** Top three panels: The positioning of the human-aligned human-engineered  $V_H$  antibodies (Kymouse; 209,452, PAD; 99,213 and therapeutic mAbs; 329) in the developability profile space of the native human  $V_H$  dataset (854,418 antibodies) based on a principal component analysis (PCA, see Methods). Bottom two panels: The positioning of the human-aligned human-engineered  $V_L$  antibodies (PAD; 78,921 and therapeutic mAbs; 320) in the developability space of the native human  $V_L$  dataset (385,633 antibodies). The hexagonal bins (shown in the back layer) represent the count of native antibodies (scale shown on the top right of the panels), and the human-engineered antibodies are represented as data points. **c** Evaluating the predictability of sequence (left panel) and structure DPs of the human-aligned human-engineered  $V_H$  antibodies (Kymouse; 209,452, PAD; 99,213 and therapeutic mAbs; 329), using multiple linear regression (MLR) models trained on native human  $V_H$  antibodies. As explained in Fig. 6a (ML Task 1), the predictive accuracy of two types of embeddings was tested, including single-DP-wise incomplete developability profiles (DPLs) and ESM-1v protein language model encoding vectors (PLM). MLR models were trained using 1000 antibodies for DPL-based predictions and 20000 antibodies for PLM-based predictions (respective saturation points). Missing DPs tested in this analysis belonged to the MWDS exclusively as determined at a Pearson correlation coefficient threshold of 0.6, for the native human IgG dataset, summing to 13 sequence DPs and 28 structure DPs (Supplementary Table 2). The y-axis represents the mean coefficient of determination ( $R^2$ ) across 20 repetitions ( $n = 20$  independent experiments). Numerical values shown represent the average values of mean  $R^2$  across (sequence or structure) DPs. Supplementary Figs. 18–21.



21.6% in mAbs) and IGH1 members (13.3% in native, 20.8 in PAD, and 25.9 in mAbs—Supplementary Fig. 18B). These findings suggest a notable resemblance between the native and human-engineered antibody datasets with regard to potential developability-related liability sequence motifs and germline annotations.

Thus, we asked how the native antibody dataset relates to the human-engineered ones in terms of DP values. To this end, we first conducted a distance-based clustering (similar to Fig. 3b) starting from the pairwise DP correlation matrices of the native and human-engineered antibody datasets (Fig. 7a). On the sequence level, we found that the general species-specific

(human and mouse) and chain-specific (light and heavy) clustering, as previously reported in our investigation of the native dataset (Fig. 3b; left panel), remained consistent when integrating native with human-engineered antibody datasets (Fig. 7a; top panel). Specifically, the human antibodies (light and heavy chains) from the PAD and mAbs datasets localized within the native human clusters of the same chain type (correlation distance  $\leq 0.1$ —Fig. 7a; top panel). Similarly, the light-chain human-engineered murine antibodies (PAD and mAbs) clustered with the IgK native mouse dataset (correlation distance  $\approx 0.1$ —Fig. 7a; top panel). Kymouse antibodies were the most distant subset, which may be explained by their unique intermediate diversity between mice and humans<sup>86</sup>.

On the structure level, we found that the original clustering pattern among the native antibody isotypes is preserved from the analysis conducted on the native-only dataset (Fig. 3b, Fig. 7a; “Native only” zone). Human-engineered subsets clustered apart from the “Native only” zone, maintaining either species-specific or chain-specific clustering. Although suggestive, the results of this analysis should be interpreted with caution due to the imbalance of dataset sizes (number of antibodies) (Supplementary Fig. 1A, B). Indeed, a correlation matrix stabilization study suggested that a dataset size of at least 50 K antibodies is required to stabilize association values among DPs (Supplementary Fig. 19A, C) and the distribution of their pairwise Pearson correlation values (Supplementary Fig. 19A, B, D). This threshold (50 K antibodies) is higher than the count of antibodies included in the murine subsets of the PAD dataset ( $\approx 24$  K for heavy chains,  $\approx 21$  K for light chains) and all the species-specific and chain-specific subsets within the mAbs dataset (329 for human heavy chains, 320 for human light chains, 62 for murine heavy chains, 71 for murine light chains) (Supplementary Fig. 1B).

Next, we leveraged the principal component analysis (PCA) that we conducted on the developability profiles of native antibodies (Fig. 5c—see Methods) to examine how human-engineered antibodies relate to the native ones in the developability spaces ( $V_H$ ; Fig. 7b: top panels,  $V_L$ ; Fig. 7b: bottom panels). In addition to comparing native and human-engineered antibodies in the *developability profile* space (dimensionality:  $10^1$ ), we also used *sequence embeddings* provided by the ESM-1v<sup>75</sup> protein language model (PLM, dimensionality:  $10^3$ ) (Supplementary Fig. 20B, Supplementary Fig. 21B). PLM embedding is an alternative to biologically motivated features (such as DPL), learned without supervision from a large pool of proteins<sup>87</sup>. Protein language modeling enables capturing longer-distance relationships within protein sequences<sup>88–90</sup>. We focused our analysis on human antibodies as they represent the largest species-specific subset among our datasets ( $\approx 0.8$  M for native  $V_H$  antibodies,  $\approx 0.4$  M for native  $V_L$  antibodies, 99213 for PAD  $V_H$  antibodies, 78921 for PAD  $V_L$  antibodies, 329 for  $V_H$  mAbs and 320  $V_L$  mAbs—Fig. 7B, Supplementary Fig. 1A, B).

Overall, we found that human-engineered antibodies (both  $V_H$  and  $V_L$ ) are majorly contained within both the developability and PLM landscapes of the native antibodies (Fig. 7b, Supplementary Fig. 20B, Supplementary Fig. 21B), suggesting that—for the DPs included in our analysis (see Methods)—the developability and sequence landscapes of human-engineered antibodies merely occupy subspaces of the natural space (in terms of the two main axes of variation studied).

From the native repertoire perspective,  $V_H$  antibodies coalesced into a single cluster in the  $V_H$  developability space, and the positioning of the antibodies in this space was independent of both their isotype and IGHV gene family annotation (Supplementary Fig. 20A). In contrast to  $V_H$  antibodies, native  $V_L$  antibodies clustered in two distinct clusters in the  $V_L$  developability space where the majority of IgK antibodies (99.999%) and IgL antibodies (95.6%) occupied the bottom cluster, and a small proportion of IgL antibodies (4.4%) predominantly occupied the top cluster (Supplementary Fig. 21A). This finding suggests that, except for a minor subset of native IgL antibodies,  $V_L$  sequences (both IgK and IgL) are homogeneous with respect to their developability. This aligns with a recent report where, using the therapeutic antibody profiler (TAP), it was found that native IgL antibodies exhibited comparable structural developability characteristics to those of native IgK antibodies<sup>36</sup>. Among the human-engineered antibodies,

only two (out of 320  $V_L$ ) therapeutic mAbs (Erenumab and Benteracimab) of the isotype IgL, and only 1345 (1342 IgL and 3 IgK, out of 78,291  $V_L$ ) patent-submitted antibodies were found to be contained in the top cluster (Fig. 7b), displaying a similar distribution trend between the two clusters as the native  $V_L$  dataset. It is worth noting that previous studies (including TAP) have previously reported Erenumab to exhibit developability risk factors<sup>7,10</sup>, which may explain its localization in the top cluster of the  $V_L$  developability space (Fig. 7b). Benteracimab was not included in these studies as it is not yet clinically approved (Phase III of clinical development as of September 2023)<sup>91</sup>.

As the human-engineered antibodies were shown to harbor comparable developability and sequence properties to those of the native ones (Fig. 7b, Supplementary Fig. 20B, Supplementary Fig. 21B), we investigated the generalizability of the native-trained multiple linear regression (MLR) models to predict the values of single missing DPs of the human-engineered datasets (Fig. 7c). Specifically, we implemented the MLR models from ML Task 1 (Fig. 6a and b) at the training sample size, which achieved the highest prediction accuracy before plateauing (1000 antibodies for DPL-based predictions, 20 K antibodies for PLM-based predictions) to predict the values of MWDS DPs (sequence and structure) on the human-engineered antibody datasets (Fig. 7c).

We found that the predictability of DP values for the human-engineered antibodies was similar to that of the native antibodies (Fig. 6b, Fig. 7c). For instance, the mean prediction accuracy (mean  $R^2$ ) of sequence DPs was between 0.25–0.33 for DPL-based predictions, and between 0.80–0.87 for PLM-based predictions among human-engineered datasets (Fig. 7c) in comparison to 0.34 and 0.92 (DPL and PLM respectively) for native antibodies (Fig. 6b). Similarly, when considering human-engineered datasets, the mean prediction accuracy of structure DPs ranged from 0.10 to 0.12 for DPL-based predictions and from 0.15 to 0.24 for PLM-based predictions (Fig. 7c). In comparison, the native antibodies exhibited prediction accuracies of 0.12 and 0.27 (respectively—Fig. 7c).

Collectively, our results suggest that the knowledge learned from the native antibodies in regards to their developability and sequence properties is, in part, generalizable to the human-engineered antibody datasets. Nevertheless, it is worth reiterating that the predictability analysis performed on the native and human-engineered datasets aims to inform experimental design and investigate factors of generalizability rather than improving predictability.

## Discussion

Previous studies have shown that several experimental DPs may be computationally inferred<sup>7,12,32,56,60</sup>, which supports the real-world relevance of computer-based developability screening of large antibody sequence libraries. However, discrepancies between computational and experimental DP profiling remain<sup>18,30</sup>. Furthermore, previous studies using computational profiling varied in the number of DPs, ranging from  $<10$  to  $>500$ <sup>10,55,57</sup>. In this study, we used 86 sequence and structure-based DPs, many of which were pairwise lowly or uncorrelated to one another (Fig. 2a and b, Supplementary Fig. 6, Supplementary Fig. 7, Supplementary Fig. 9), thus capturing at least a subspace of the multidimensional developability space. So far, the relevant dimensionality of the antibody developability space remains unclear. That said, our analysis can be replicated with any other set of computationally or experimentally determined DPs. While most previous studies focus on small antibody datasets and mostly on the comparison between experimental and computational DPs, this study explores the plasticity of the *in silico* DP space on large-scale antibody datasets. While our findings may partly depend on the DPs studied, our conclusions derived from DP profiling (Fig. 7b) are consistent with those from PLM-based profiling (Supplementary Fig. 20B, Supplementary Fig. 21B). Specifically, we observe that in both representations, human-engineered antibodies are predominantly localized within regions occupied by natural antibodies and display a tendency to cluster in specific areas rather than dispersing uniformly throughout the space. In other words, our study should be understood as forecasting the types of analyses possible once large-scale

experimental or highly validated computational DPs exist, allowing for routine repertoire-scale DP profiling and prediction. In summary, the real-world relevance of our systems approach is grounded in (i) the profiling of a large number of pairwise-independent DP parameters, (ii) sequence and rigid and dynamic structure-based DP analysis (Supplementary Note 2), (iii) diverse experimental datasets ranging from native to human-engineered, and (iv) parameter-independent computational and machine learning analysis. In the future, it would be of interest to add display libraries for comparison<sup>92–94</sup> as well as a larger number of paired datasets with broad and deep isotype information<sup>95,96</sup>.

Intercorrelation analysis revealed that structure-based DPs show higher independence (lower pairwise correlation) than sequence-based DPs (Fig. 2). These findings emphasize the significance of structural consideration for therapeutic antibody design<sup>3,8,25,97,98</sup>. Due to the scarcity of antibody structural information, structure-based developability estimation has previously presented a substantial challenge in developability screening<sup>3,20</sup>. However, *in silico* high-throughput antibody structure prediction tools have evolved in speed and accuracy with machine learning algorithms, permitting the screening of structural parameters in large datasets<sup>9,38,39,99</sup>.

The scale of this work's analysis facilitated the discovery of parameter redundancies that could accelerate future antibody developability screening processes. For instance, while Chen and colleagues included both the molar extinction coefficient and the extinction coefficient of the variable region sequence (AbChain\_moltextcoef, AbChain\_percenttextcoef) and the cysteine bridges (AbChain\_cysbridges\_moltextcoef, AbChain\_cysbridges\_percenttextcoef) as important developability predictors<sup>57</sup>, we found that one of these coefficients could likely be sufficient to replace the other. Similarly, although Ahmed and colleagues highlighted the importance of structure-based isoelectric point (pI) as an essential developability parameter on a limited-size therapeutic antibody dataset (77 clinical-stage antibodies)<sup>10</sup>, our analysis suggested that sequence-based pI (AbChain\_pI) could potentially replace structure-based pI. We further highlighted the importance of large-scale antibody developability data to stabilize the associations among DPs (Supplementary Fig. 19), and, thus, to reveal such redundancies. Additionally, most developability studies emphasize the relevance of antibody charge in the physiological pH (7.4)<sup>7</sup>. However, mAbs are usually exposed to a wider range of solution pH (4.8–9) during production and formulation<sup>10,100</sup>. Also, antibody variable region charge in acidic pH has proven to be a critical factor in IgG mAb pharmacokinetics<sup>101–103</sup> and product formulation<sup>28</sup>. Therefore, we included the sequence-based variable region charge measures in 14 pH points (1–14) in our analysis. We found that the charge measures of native IgG sequences formed three correlation clusters with intermediate (0.4–0.7) pairwise correlations, emphasizing the importance of antibody charge considerations on a wider spectrum of pH values (Fig. 2b). Other pairs of potential parameter redundancies are the sequence content of polar and non-polar amino acids (AbChain\_polar\_content Vs AbChain\_nonpolar\_content), the sequence content of aliphatic amino acids and the sequence aliphatic index (AbChain\_aliphatic\_content Vs AbChain\_aliphatic\_index), the average atomic interaction distance of the antibody structure and the number of Van der Waals clashes (AbStruc\_mean\_interaction\_distance Vs AbStruc\_vdw\_clashes) and the sequence length and the sequence molecular weight (AbChain\_mw Vs AbChain\_length) (Fig. 2A). Redundancy-based parameter reduction, analogous to feature selection for ML models, would accelerate future antibody developability investigations by screening for a more comprehensive and sufficiently representative set of parameters.

Our analyses revealed chain-specific developability signatures in relation to DP values and pairwise associations (Fig. 3b, c), emphasizing possible differences in developability design considerations for therapeutic antibody development. Within each chain type, we found that murine and human antibodies occupy distinguishable developability spaces (Fig. 3c), highlighting the importance of transgenic mice for antibody screening and the challenges of antibody humanization efforts<sup>104–107</sup>. Interestingly, our findings suggest that the Kymouse (humanized mice) dataset under investigation undersampled the human dataset, both with respect to developability

(Fig. 7b) and sequence spaces (Supplementary Fig. 20B), even though lower-level features such as VDJ gene usage and CDR3 length were previously found to overlap<sup>86</sup>.

In addition to chain type and species, we found that human heavy chain (V<sub>H</sub>) isotypes harbor high similarities in regard to their pairwise DP associations (Fig. 3b) and redundancies (Supplementary Fig. 10B, Fig. 3a). We found that they aggregated homogeneously in the developability space regardless of their isotype (Supplementary Fig. 20A). Thus, although all currently approved therapeutic mAbs belong to the IgG isotype<sup>108</sup>, our findings provide the incentive to explore the available native antibody Fv sequence space beyond the isotype annotation for novel mAb discovery. In regards to V<sub>L</sub> isotypes, human IgK and the majority of IgL antibodies clustered together in the developability space (Supplementary Fig. 21A), questioning the previous association of IgL antibodies with poor developability, in line with recent findings by Raybould and colleagues<sup>36</sup>, and providing an incentive to re-include IgL sequences in future antibody discovery libraries. Nevertheless, we are aware that our findings involve only the Fv regions of antibody sequences as current antibody-specific structure prediction models do not take into account the Fc region and do not account for the impact of the Fc region on developability<sup>103,109</sup>.

When examining the impact of sequence similarity on developability similarity, we found that antibodies that are highly similar in sequence can possess dissimilar developability profiles (Fig. 5), which is in line with previous findings<sup>110</sup>. This suggests that there are degrees of freedom available for therapeutic antibody candidate engineering to optimize developability with minimal changes in antibody sequence.

For instance, if an antibody candidate exhibited optimal antigen binding properties with a suboptimal developability profile, minor sequence changes could result in a substantial improvement (or deterioration) of this profile without major changes to its antigen binding properties. In this context, Petersen and colleagues showed that the native human antibody repertoire can aid the identification of advantageous (or disadvantageous) “universal” (native) framework mutations that could facilitate therapeutic mAb development<sup>54</sup>. Furthermore, small changes in antibody sequence with large effects on function have been observed for both developability and antibody binding properties<sup>3,111</sup>, underscoring the importance of simultaneous optimization of multiple design parameters. Since optimization of a specific property often comes with undesirable trade-offs for other properties<sup>29,61</sup>, studying the effect of substitutions can help with designing antibodies that exhibit improved developability with comparable efficacy.

To directly probe the relationship of sequence to developability, we performed a single amino acid substitution analysis (see Methods). Although changing one amino acid at a time only explores a small fraction of the input factor space<sup>62</sup>, it allows for the definitive attribution of the change in output to a specific amino acid substitution. Although including all possible variants with two or more substituted amino acids could account for epistatic mutation effects, it would geometrically inflate the combinatorial space beyond current computational feasibility. Because exhaustively mapping the DPs of such large sequence spaces is (currently) computationally infeasible, to explore more than single amino acid substitutions, one must find ways to efficiently and uniformly sample the input space, possibly through latin hypercube sampling or low-discrepancy sequences<sup>112,113</sup>. Briefly, we found that across all possible single amino acid substituted variants of a sample of 500 human heavy chain antibodies some parameters were especially insensitive to the average mutation. Furthermore, to quantify the sensitivity of a parameter, we measured the dispersion of parameter values of all possible single amino acid substituted variants of an antibody. We used “tailedness” (measured by excess kurtosis) and the range of the distributions as proxy measures for average and maximum potential sensitivity, respectively. Since excess kurtosis is, strictly speaking, defined for normal distributions, its interpretive power depends on the nature of the observed parameter. It may be beneficial in future work to consider additional properties of distributions, such as skewness or diversity (as measured by Shannon entropy). Despite the apparent low sensitivity of charge and hydrophobicity DPs, they have been shown to impact the developability of

monoclonal antibodies greatly<sup>100</sup>. This is, in fact, not contradictory since the “tailedness” does not describe absolute DP shifts, but the proportion of outliers of a DP distribution (average sensitivity) or the range proportional to the original wildtype (WT) sequences (potential sensitivity). Since the sensitivity metrics are an aggregate score averaged over all variants of the sampled antibodies, they serve as a rough global characterization of DPs. Drastic property changes still require more comprehensive sequence- and structure-based local approaches in individual antibodies to model<sup>114</sup>. Prospectively, although we only studied general trends in sensitivity from the perspective of individual developability parameters, it would be interesting to explore how a mutation impacts values across all DP of an individual antibody.

In this work, we did not perform structure-based substitution analysis. Specifically, given that there is a lack of experimentally determined structures of antibody variants, it remains unclear how accurately antibody structure prediction tools can resolve single amino acid structural differences and reflect them in models (Supplementary Fig. 11, Supplementary Note 3)<sup>115</sup>.

More generally, in the context of classical MD simulations of antibodies, the conformational landscape typically exhibits minor variations compared to the initial rigid structure<sup>116</sup>. In MD simulations of five experimentally determined antibodies without antigen (Supplementary Note 2), we found that the structural fluctuations in the CDRH3 loop regions were generally minimal (on the order of 0.1 Å) throughout 100 ns simulations (as determined by RMSD)<sup>117</sup>. However, the structural variance between the initial (rigid) antibody conformation and the relaxed version was between 1.2 Å and 1.7 Å on average (Supplementary Fig. 4A, Supplementary Note 2). Given the relatively small magnitude of CDRH3 fluctuations and the current limitations of structure prediction methods, it is challenging to detect and accurately represent these small conformational changes in rigid antibody models. The accuracy of current structure prediction tools typically is below the range of these small fluctuations<sup>118</sup>. However, when it comes to studies involving antibody variants, we cannot overlook these fluctuations, as they reflect the effects of mutations. These computational structure analyses should be combined with experimentally determined antibody structures, such as X-ray or cryo-EM data. While experimental structures contain variations due to differences in crystallization conditions, resolution levels, and inherent protein dynamics, making the same structure slightly different in various studies, computationally predicted structures lack the experimental-related noise found in real structures, causing all models to appear identical outside the mutation site. Given these limitations, we anticipate that the differences in distances between variants, specifically at the mutation site, will be lower than expected in experimental variants and greater than expected in computational models. Our findings suggest that the reality of variant effects lies between the generated models and structures observed experimentally. Furthermore, due to the dynamic nature of proteins, simulated or experimentally determined structural differences of the same protein can overshadow differences caused by mutations, especially in highly flexible antibody loops like CDRH3.

ML models have been shown to be able to predict the biophysical characteristics of proteins<sup>29,31,119,120</sup>. While their use may not be widespread in the pharmaceutical industry, the emergence of AI/ML may become routine as part of initial *in silico* efforts to screen and assess molecular properties and interactions before any experimental efforts<sup>23</sup>. To efficiently determine biological (and, possibly, clinical) properties of antibodies using machine learning-based methods, it is important first to identify a suitable representation. In our study, we compare two alternative embedding types, one (DPL), which collects (a subset of) DPs into a numerical vector and the other (PLM) obtained by encoding the antibody sequence using a pre-trained neural network (NNs). The former representation mirrors feature-selection approaches to DP prediction, while the latter embraces the latest advances in deep NNs and, in particular, transformer-based models. We compute DP predictions by passing DPL and PLM embedding through *linear regressor heads*. As shown in Fig. 6b, the predictive power of our models is limited and

this holds especially for DPL predictors, even on sequence DPs, which are instead well predicted using the PLM representation. The stark difference in performance is not surprising and can be understood as the combination of two facts. First, we observe that our use of linear heads is a restrictive architectural choice, which is likely over-simplistic: it is natural to expect that most—even though not all—DPs cannot be written as a linear combination of the remaining ones. This limitation is exacerbated by multiple collinear features (as in the full developability profile, Fig. 2) and by a reduced number of dimensions (the number of coordinates of DPLs vs PLMs span two orders of magnitude). Second, as visible in Supplementary Fig. 20C (right), PLM representations retain considerable information about sequence identity and they can be thought of as a learned map from sequences to a latent space. They are, therefore, clearly more suitable to derive sequence DPs and, as our experiments find, also structure DPs, although to a lesser extent. Overall, PLM-based DP predictions generalize well across human-engineered datasets, hinting at a modest amount of overfitting. This, paired with the observation that regressor quality improves with bigger train set sizes (no clear plateaus found for the cardinalities tested in Fig. 6b), indicates a healthy learning trend and suggests that predictions may benefit from more advanced head architectures.

DP predictors also shed light on the relationship between different groups of DPs, e.g., between sequence and structure DPs, and MWDS vs. non-MWDS ones. Both DPL- and PLM-based regressors fail in inferring most structure-based parameters and, with few exceptions (notably `AbStruc_psi_angle` and `AbStruc_unfolded_pI`), succeed on the same DPs (e.g., `AbStruc_loops`, `AbStruc_beta_strands`, `AbStruc_sasa`, `AbStruc_unfolded_pI`, `AbStruc_pcharge_hetrgen`, `AbStruc_psi_angle`). This seems to suggest the failure of both representations to capture the complex relations between sequences in linear subspaces of the latent space. We also remark that although the lower correlation between structure DPs seen in Fig. 2A may suggest that noise introduced by structure prediction tools also contributes to the weak prediction outcomes, we were unable to obtain better scores when predicting structure DPs computed on the (few) available solved antibody structures.

Of note, the ML methods employed in this work are by no means exhaustive and only represent a first step toward ML-based DP predictability. It is of interest that already relatively simple ML approaches achieve good prediction accuracy. More generally, our ML approach is to be understood as an example of potential studies that may be performed once it is clearer which DPs are causally linked to downstream antibody candidate success.

When examining the role of structure prediction tool on DP value computation, we reported that structure-based DPs calculated on different antibody structure prediction models correlate overall poorly (Supplementary Fig. 3A, Supplementary Note 2). These observations align with recent reports by others<sup>32,36,118</sup>. Importantly, we found that although `ABB2`<sup>38</sup> and `IgFold`<sup>39</sup> have been found to be superior to `ABB`<sup>37,38,121</sup> (which was used in this work for the majority of the results), the correlation of `ABB` with all other tools (experimental) in terms of structure-based DPs was low and did not differ to the aforementioned tools. Strikingly, we found that structures obtained by computational structure prediction methods belonged to a common underlying ensemble structural distribution as revealed by MD (Supplementary Fig. 4, Supplementary Note 2). Therefore, the usage of `ABB` did not considerably bias the results in this paper.

However, it remains challenging to predict the antibody CDRH3 loop regions, which are of great interest due to their involvement in antibody-antigen binding<sup>97</sup> as well as developability<sup>20</sup>. These loop regions are generally more flexible and less structurally conserved compared to stable secondary structure elements like beta sheets<sup>118</sup>. Consequently, predicting the specific conformations of these loops accurately can be difficult, and it becomes even more challenging to determine when a prediction is correct without combining experimental validation and MD simulations<sup>32,118,122</sup>.

MD is important for a fuller understanding of antibody systems as it provides insight into the flexibility and fluctuations of antibody structure and developability parameters<sup>28,32,36,121,123,124</sup>. Specifically, Park and Izadi

found that antibody developability surface descriptor parameters (e.g., positive electrostatic potential on the surface of the CDR region<sup>125</sup>) vary extensively as a function of the structure prediction method used, which is in line with the findings in this manuscript (Supplementary Note 2). However, after Gaussian-accelerated MD (GaMD) simulations and averaging the values of the descriptors through MD frames, the consistency between the descriptors improved<sup>32</sup>. Similarly, Raybould and colleagues<sup>36</sup> evaluated variations in four structure-based Therapeutic Antibody Profiler (TAP)<sup>7</sup> scores using molecular dynamics (MD). The mean values of the TAP properties observed during the simulations closely matched an ensemble generated from three TAP predictions based on static Fv models directly generated by ABodyBuilder2. Furthermore, certain structure prediction tools, such as IgFold<sup>39</sup>, refine the final model using short relaxation using OpenMM<sup>126</sup> or PyRosetta<sup>127</sup>.

In the future, generating multiple models with various refinements could yield a conformational ensemble similar to what is obtained from a full-fledged long MD simulation<sup>118,122</sup>. Training on dynamic data, such as MD simulation-based conformations, enables the model to capture loop flexibility and variability, resulting in more robust and realistic predictions that can improve currently challenging CDRH3 structure prediction. An ideal ML model trained on dynamic data could streamline antibody structure prediction by eliminating the need for additional MD simulations, saving computational resources, and providing a higher level of confidence in predicted structures and insights into hidden antibody characteristics and developability properties.

Since the development of therapeutic antibodies involves human-directed sequence engineering to attain desirable developability properties, there have been efforts to separate natural from therapeutic antibodies<sup>7,55</sup>. However, we showed that human-engineered (including therapeutic) antibodies fall *within* the developability space of natural antibodies (Fig. 7b). Although this could simply be due to the particular selection of DPs and the (major) principal components of variations, we showed that human-engineered antibodies also fall within the *sequence* space of natural antibodies (Supplementary Fig. 21B). Here, we discuss several possible explanations:

(1) Our findings are consistent with<sup>52</sup>, who have shown considerable sequence overlap between therapeutic antibodies and NGS-derived natural repertoires, indicating that therapeutic antibody sequences are largely derived from natural sequences with few modifications. This could also explain why MLR models trained on natural DPs predict DPs of human-engineered antibodies with similar accuracy (Fig. 7c). (2) Although the development and formulation of therapeutic antibodies involve distinct challenges from natural antibody generation, both might be subject to converging primary restrictions. As shown by ref. 54, certain framework mutations regularly occurring in clinical-stage therapeutic antibodies are also frequently observed in natural repertoires, indicating converging selection towards common characteristics such as stability. It is conceivable that the sequence space of stable (and non-immunogenic) antibodies is restricted so that natural and human-engineered sequences occupy overlapping subspaces within. (3) There might be considerable engineering of the Fc region<sup>128</sup>, which is not included in our study, that could separate human-engineered from natural antibodies in sequence space. (4) Lastly, the number of patent-submitted<sup>58,129</sup> and therapeutic antibody sequences available is relatively small, which means that there may be potential therapeutics in yet unexplored regions of the sequence space. However, given that sample sizes differ across the human-engineered datasets, more data is needed to study how cross-transferable conclusions are.

Despite being globally similar to the native dataset, DPLs of human-engineered antibodies show distinctive traits. Most notably, they differ in the patterns of correlation found among DPs. This is best seen in Fig. 7b, which is obtained by projecting native DPLs along two axes which are determined in order to be uncorrelated. In this subspace, the cluster originated by native antibodies has a circular shape, i.e., it is isotropic, meaning that the knowledge of one of the two coordinates gives limited information about the other. This property, however, does not hold for the set of engineered

antibodies, which form an elongated shape, stretching along a diagonal. A similar cluster shape implies that growing x-values tend to correspond to bigger y-values or, in other terms, that there is a (weak) positive correlation between the two axes. Among human-engineered antibodies, then, alterations to the two parts of DPL, which are projected on each axis, do not happen independently anymore, likely signaling the artificial effect on DPs of the antibody optimization process. In the future, it would be of interest to study how this potential signal of the antibody optimization process is represented in different PLM embeddings, for example, in antibody-specific PLM embeddings, for which evidence is conflicting as to whether general protein or antibody-specific PLM more faithfully represent inter-sequence functional similarity<sup>77,130,131</sup>.

Furthermore, it is interesting to note that human-engineered DPLs do not cover entirely the native DPL cluster (Fig. 7b). Among human-engineered antibodies, it is hence less probable to find some types of DPLs, which appear instead frequently among native ones: this suggests the existence of an unexplored region of DP space. Our conclusion is strengthened by the presence of a similar low-density zone on the right of the cluster of PLM embedding (Supplementary Fig. 20B), which may constitute evidence of under-investigated classes of antibody sequences.

Finally, computational developability profiling, similar to computational antigen binding profiling<sup>3,98</sup>, will be increasingly useful once the real-world relevance of computational DPs has been further established<sup>30–32,100</sup>. To this end, multiple avenues require further development: (i) insight into differences of paired vs single-chain developability<sup>95,96</sup>, (ii) faster computation of structural and MD-based DPs<sup>132</sup>, (iii) negative controls for clinical stage antibodies that have not progressed due to developability issues, (iv) large-scale data where multiple properties for a given antibody are captured to start exploring multiparameter generative design<sup>29,61,106,133</sup>, and (v) open and unbiased competitions with agreed-upon quality experimental (and simulated<sup>69,134</sup>) data to quantify the current state-of-the-art in DP predictability, causal relationships and importance ranking vis-a-vis DP impact on antibody developability.

## Methods

### Antibody datasets

**The native antibody dataset.** We collected 2,036,789 native human and murine antibody sequences (Supplementary Fig. 1A). Briefly, we assembled non-redundant sequences of human (heavy chains: IgD 173,342, IgM 173,437, IgG 170,473, IgA 171,174, IgE 165,992, and light chains: IgK 198,255, IgL 187,378) and murine (heavy chains: IgM 198,967, IgG 199,326, and light chains: IgK 198,795, IgL 199,650) native antibody variable region sequences (Supplementary Fig. 1A). Antibody sequences were majorly sourced from Observed Antibody Space (1,738,091 sequences)<sup>135</sup> in addition to including our own experimentally-generated sequences (total of 298,698 sequences with the IgD, IgK and IgL human datasets) to provide balanced antibody count among all isotypes. Experimental sequences were generated following a protocol inspired by ref. 136, starting from human blood (see Methods), and exported as VDJ clones from raw sequencing files using MiXCR (version 3.0.1)<sup>137</sup>. OAS sequences were aligned and IMGT-numbered using the IMGT/High V-QUEST<sup>138</sup>. The kappa and lambda light chain types (IgK and IgL) were often referred to as isotypes in this manuscript to provide a cohesive reading experience.

### The human-engineered antibody datasets

(1) The therapeutic antibody dataset (mAb). A total of 782 therapeutic antibody sequences belonging to the “whole mAb” format were obtained from TheraSABDAB as per July 2021<sup>139</sup>. First, all mAb sequences were aligned to both murine and human germlines (species of interest) with the antigen receptor alignment and annotation tool ANARCI<sup>140</sup>. We used the same IMGT-numbering scheme as implemented in the native dataset sequences numbering to ensure consistent methodology. For each chain, we kept the alignment with the highest germline certainty (human or mouse) by selecting for the highest V gene identity (the sequence identity over the



V-region to the most sequence-identical germline). In case V gene identity was equal for murine and human alignments, we kept the alignment with the highest J gene identity (the sequence identity over the J-region to the most sequence-identical germline). To ensure the accuracy of annotations, we excluded sequences with V gene identities <0.7 (Supplementary Fig. 1B).

(2) The patented antibody database (PAD). Under a non-commercial agreement with NaturalAntibody, we obtained a total of  $\approx$  240 K unpaired variable region antibody sequences (Supplementary Fig. 1B). Sequences were originally extracted from patent documents from intellectual property organizations and bioinformatic databases (WIPO: world intellectual property organization, USPTO: United States Patent and Trademark Office, EBI: European Bioinformatics Institute, DDBJ: DNA Data Bank of Japan), and mined as explained by ref. 58. Sequences were provided with species and V-gene annotation metadata and IMGT-numbered. We selected for human and murine sequences and classified them into heavy (IgH) and light (IgK and IgL) chains based on the corresponding V-gene annotation resulting in a final count of 223,613 PAD sequences (Supplementary Fig. 1B).

(3) Humanized mouse dataset (Kymouse). We obtained 209,452 IgM  $V_H$  sequences from humanized transgenic mice (Kymouse)<sup>86</sup>. Sequences were downloaded from OAS<sup>135</sup>, and their structures were predicted with ABB as explained above. Finally, we also measured their sequence and structure DPs.

**The in silico mutated antibody dataset.** We first sampled 500 wildtype (WT) antibodies from the human native  $V_H$  dataset (100 samples per isotype for IgM, IgD, IgG, IgA and IgE). We generated all possible acid substitutions for each antibody, resulting in a total of 301,777 sequences for which we predicted/computed their sequence-based developability parameters as detailed above. We used this data to perform the sequence DP sensitivity analysis described in Fig. 4.

### Antibody structure prediction

We used ABodyBuilder (ABB) to predict the structures of antibody variable regions<sup>37</sup>. The high-throughput version of ABB is provided as an image for Vagrant VirtualBox (version 2.2.16), known as SabBox. We ran this version with default parameters using unpaired single chain input (heavy or light) to predict all structures in all datasets, unless mentioned otherwise. The ABB pipeline includes a relaxation with MODELLER after the initial antibody structure prediction.

### In silico calculation of developability parameters

**Calculation of sequence-based developability parameters.** *Molecular parameters:* the molecular weight of an antibody variable region sequence was calculated using the Peptides R package version 2.4.4<sup>141,142</sup>. We calculated the antibody length and the average residue weight using custom R scripts<sup>143</sup>. *Amino acid categorical composition:* Using the Peptides R package and a custom R script, we calculated the proportional (%) content of amino acid categories<sup>142</sup> (Supplementary Table 1) by dividing the occurrences of the amino acids in one group by the length of the antibody variable region sequence. *pI and charge:* To compute the pI and charge of a given sequence, we used the Peptides R package<sup>141</sup> using the “Lehninger” scale between pH = 1 and pH = 14 with step size = 1 (14 data points). *Extinction coefficient and molar extinction coefficient (for all and for cysteine bridges):* Calculations were performed as mentioned in refs. 144,145 using custom R scripts. *Hydrophobicity and hydrophobic moment:* To compute hydrophobicity and the hydrophobic moment, we used the Peptides R package<sup>142</sup>, the scale of choice was Eisenberg. For hydrophobic moment calculation, we selected ten amino acids for the length of the sliding window size based on our understanding of the antibody secondary structure and as explained by ref. 146. We also specified the angle value as 160 as recommended by ref. 147. *Instability index and aliphatic index:* The aliphatic index is defined as the relative volume occupied by aliphatic side chains (Alanine, Valine,

Isoleucine, and Leucine—Supplementary Table 1). It may be regarded as a positive factor for the increase of thermostability of globular proteins<sup>148</sup>. The instability index was first developed by ref. 149 to reflect the stability of proteins based on their content of certain dipeptides that were found to be associated with degradation tendency. We calculated the instability and aliphatic indices using the Peptides R package<sup>141</sup>. **Protein solubility prediction:** We used SoluProt (version 1.0) to predict the solubility index for antibody variable region sequences<sup>150</sup>. We chose this tool as it allows the consideration of protein expressibility into its solubility score prediction while proving comparable performance compared to other state-of-the-art solubility prediction tools<sup>150</sup>. Sequences with solubility scores above 0.5 were predicted to be soluble when expressed in *Escherichia Coli*. **Immunogenicity prediction:** We used netMHCIIpan version 4.0<sup>40</sup> to estimate the immunogenicity of the antibody variable region sequences. Briefly, we examined the global immunogenicity of the antibody sequences by predicting their affinities for HLA II supertypes that are found in 98% of the global population<sup>111</sup>. We obtained the following numerical values from these calculations including (i) minimum rank, (ii) number of weak binders (percentage rank >2 and <10), (iii) number of strong binders (percentage rank < 2), (iv) full average percentage rank and (v) the number of antibody regions where the maximal immunogenic peptide stretches (maximal\_immunogenicity\_region\_span). As these calculations are computationally intensive, we computed the immunogenicity DPs on 10% only of the native, Kymouse and PAD datasets.

**Calculation of structure-based developability parameters.** First, we added hydrogen atoms to each antibody variable region structure that we previously built with ABodyBuilder using the Reduce software (version 3.24.130724)<sup>151</sup>. Hydrogenated structures were used as an input (pdb format) to calculate structure-based developability parameters as they were reported to provide a closer chemical representation of their potential in vivo structures<sup>152</sup>. As detailed in Supplementary Data 1, we used BioPython (version 1.79) to compute secondary structure parameters<sup>153</sup>, FreeSASA (version 2.1.0) for solvent accessibility predictions<sup>154</sup>, PROPKA (version 3.4.0) for the calculation of thermodynamic and electrochemical parameters<sup>155</sup> and ProDy (version 2.0) to count free and bridged cysteines<sup>156</sup>. We used custom Python scripts for the calculation of developability index (DI) and spatial aggregation propensity (SAP) following the work of Lauer and colleagues<sup>56,157</sup>. We used the original Black-Mold hydrophobicity scale as suggested by ref. 124. Finally, we used Arpeggio (version 1.4.1) to calculate interatomic interactions after converting antibody hydrogenated structures to cif format<sup>158</sup>.

### Parameter correlation calculation and visualization

For the analysis in Fig. 2b, Supplementary Fig. 9A, Supplementary Fig. 6 and Supplementary Fig. 7, we computed pairwise developability parameter correlations (Pearson) matrices using the ‘cor()’ function from the ‘base’ package in R (version 4.0.3). for each isotype/species combination to investigate parameter associations and redundancies. Missing data was accounted for by choosing the argument *use.pairwise.obs = complete* option in the function. We visualized the correlation matrices using ComplexHeatmap R package version 2.9.4<sup>159</sup> and annotated the heatmaps with the threshold-specific ABC-EDA/MWDS algorithm output<sup>160</sup>.

### Determination of the minimum weight dominating set of developability parameters

To investigate the redundancy of developability parameters (DPs), we first constructed undirected weighted network graphs for each species and isotype where the nodes represent DPs and the edge weights represent the pairwise Pearson correlation values (as in Supplementary Fig. 9B). For this purpose, we constructed and visualized networks with Cytoscape 3.9.1<sup>161,162</sup>. The constructed networks could include up to three distinct classes of pairwise relationships for a given correlation threshold (from 0.1 and 0.9 in intervals of 0.1): (1) isolated nodes representing DPs that have correlation

values below the given threshold with all other DPs, (2) doublets (exclusive pairs) where two DPs are solely correlated with each other (above the given threshold), and (3) correlation subnetworks where more than two DPs form correlation clusters.

Secondly, the minimum weight dominating set (MWDS) for a specific

$$\text{Normalized Levenshtein distance} = \frac{LD(A, B)}{\text{Max}(\text{length}(A), \text{length}(B))} \quad \text{Sequence similarity} = 1 - \text{Normalised Levenshtein distance}$$

network at a given correlation threshold was defined as the set of parameters for which the sum of the associated edge weights is minimal<sup>59</sup>. Thus, the MWDS is a subset of nodes where each node in the network is either in the MWDS or connected to a member of it by an edge<sup>59</sup>. Based on this definition, we included all DPs from classes (1) and (2) in the MWDS, where no meaningful selection based only on correlation can be made. However, finding the dominating parameters among class (3) DPs is an NP-hard problem. Thus, we implemented an algorithm described by Shetgaonkar and Singh, which leverages the local optimization of an artificial bee colony algorithm guided by iterative global estimations of distribution (ABC-EDA algorithm) to approximate optimal solutions<sup>59</sup>. Ultimately, this algorithm classifies class 3 DPs as either “dominant” or “redundant”. The dominant DPs were added to the final MWDS.

### Correlation distance dendrograms

This analysis aims to examine the similarities in the pairwise associations of developability parameters among isotypes and species of the native and human-engineered antibody datasets (Fig. 3b and Fig. 7a). For this analysis, we transformed the isotype- and species-specific parameter correlation matrices (sequence and structure parameters separately) into numerical vectors. Subsequently, we quantified the pairwise Pearson correlation distance for these numerical vectors using the ‘get\_dist’ function from the R package factextra version 1.0.7<sup>163</sup>. We finally clustered the resulting distance matrices using the hierarchical clustering ‘hclust’ command following the complete linkage method from the stats package, version 4.0.3<sup>164</sup>, where the height of the dendrogram represents the Pearson distance (0–1).

### Analysis of parameter sensitivity to single amino acid substitution: excess kurtosis and range

**Excess kurtosis:** To investigate the impact of changes in the amino acid sequence on DP values (sensitivity—analysis in Fig. 4 and Supplementary Fig. 12), we normalized (mean-centered) the DP values of the mutated antibody dataset by subtracting their mean and then dividing by their standard deviation. The DP distributions of all mutants of each wt-antibody were analyzed by calculating and comparing two metrics. First, the excess kurtosis, defined as  $\text{Excess kurtosis} = \frac{\mu_4}{\sigma^4} - 3$ , where  $\mu_4$  is the fourth central moment and  $\sigma$  the standard deviation of the normalized DP values. An excess kurtosis of 0 represents that of a normal distribution. It increases with peakedness and decreases with uniformity of the distribution, under the assumption that it is bell-shaped. Thus, while a high excess kurtosis indicates a small change induced by single amino acid substitutions on average or low average sensitivity, low excess kurtosis suggests high average sensitivity. The second metric is the *range* of a distribution, defined as the distance between the highest and lowest DP values and represents the maximum normalized shift that is inducible by a single amino acid substitution.

### Pairwise developability profile correlations and pairwise sequence similarity studies

We define the *antibody developability profile (DPL)* as a numerical vector that carries the values of developability parameters in a fixed order for a given antibody sequence. Developability parameter values were mean-centered and scaled to unit variance (normalized)<sup>165</sup>. Of note, scaling for a given group of antibody sequences was performed after accounting for chain type and species.

For the analysis in Fig. 5a, b, we defined the *pairwise developability profile correlation (DPC)* as the Pearson correlation value between a pair of sequence developability profiles. We computed the amino acid-based sequence similarity between two antibody variable region sequences (A and B) as follows:

Where LD(A,B) represents the Levenshtein distance between A and B. We used the stringdist R package version 0.9.8<sup>166</sup> and the Levenshtein python package version 0.20.8 to calculate LD<sup>167</sup>.

In Fig. 5c we inspected the relationship between sequence similarity and normalized developability profile (DPL) similarity by first examining the proximity of human heavy chain antibodies that belong to the same sequence similarity cluster on the 2D PCA projection plane of the developability space ( $R^N$ ), where  $N$  is the number of DPs that makes a developability profile. In this analysis, we used the MWDS DPs (which have full values for all antibodies) as identified by the ABC-EDA algorithm for the human IgG dataset at an absolute Pearson correlation coefficient threshold of 0.6, accounting for 46 DPs ( $N = 43$ , Supplementary Table 2). The PCA was computed using the Python packages scikit-learn version 1.1<sup>168</sup> and dask-ml version 2022.5.27<sup>169</sup>. Sequence similarity groups were identified using USEARCH version 11.0<sup>170</sup> as groups of 10000 (or more) antibodies that share at least 0.75 sequence similarity as defined by Levenshtein distance (Supplementary Fig. 16B).

Then, to quantify the relationship between sequence similarity and developability profile similarity, we studied the correlation between the pairwise Euclidean distance (ED) in the developability space ( $R^N$ ) and the pairwise normalized LD for a sample of 5000 human IgM antibodies as in Supplementary Fig. 16C (computationally intensive process; 12,500,000 data points for each ED and LD calculation). We ensured that these sampled antibodies belong to the same IGHV gene family to exclude the effect of IGHV family variance on the LD computations<sup>64,65</sup>. Antibody pairs with  $ED \geq 15$  were considered outliers and were excluded from this analysis (forming 0.8% of total data points).

Of note, the PCAs from this analysis were used to examine the positioning of the human human-engineered antibody datasets within the developability space of the human native antibodies (Fig. 7b) and the role of native antibody isotype and germline gene annotation in the positioning within the developability space (Supplementary Fig. 20A, Supplementary Fig. 21A).

### Predictability of developability parameters

We used the human  $V_H$  antibody sequences (854,418 antibodies) and their computed DP values (normalized; mean-centered) to assess the predictability of developability parameters in two machine-learning tasks (ML task 1 and ML task 2). Of note, we excluded mouse antibodies and  $V_L$  sequences to avoid species- and chain-specific biases, ensuring greater data homogeneity. In both tasks, the predictability of DPs was assessed by computing the coefficient of determination ( $R^2$ ) between observed and predicted DP values<sup>171</sup>.

**ML task 1; predicting the values of the same (single) missing DP.** For this task (Fig. 6a, b, and Supplementary Fig. 17A), we randomly split the human  $V_H$  antibody sequences into two subsets: (i) *training set*; containing ~80% of sequences (683,534), which was further subsampled to derive training sets of variable sizes (50, 100, 500, 1000, 10000, 20000). Specifically, for each size, we defined 20 independent training subsets (ii) *test set*; containing the remaining (~20%) sequences (170,884). Then, we used two types of embeddings to train multiple linear regression (MLR) models: (1) single-DP-wise incomplete developability profiles (low-dimensionality embeddings—order of 10 s) and (2) antibody sequence

encodings obtained from the protein language model (PLM) ESM-1v (high dimensionality embeddings—order of  $1000 s^{74-76}$ ). Beginning-of-sentence (BOS) processing was used to compress the vectors produced by the PLM to avoid biases that might occur by averaging the entire vectors (as detailed below in “Antibody sequence encoding with PLM”). Finally, we compared the predictive power of both embeddings to predict the values of DPs in the test set after deleting single-DP column values at a time. For each type of embedding, we predicted the missing DP using linear regressors, trained with the MSE loss. We did not opt for more complex types of regressors to contrast the occurrence of overfitting. The DP value prediction was repeated 20 times (using the 20 independently trained MLR models per training set size) and the mean  $R^2$  was reported.

The native-trained MLR models from this task were then utilized to predict the values of MWDS DPs (Supplementary Table 2) in the human-engineered datasets (shown in Fig. 7c). Models were trained using the training set size, which achieved the highest prediction accuracy before plateauing (1000 antibodies for DPL-based predictions 20 K antibodies for PLM-based predictions).

**ML task 2; predicting the values of randomly missing DPs.** For this task (Fig. 6a, c and Supplementary Fig. 17B), we randomly deleted (either 2% or 4% of) DP values from subsamples of the human  $V_H$  antibody sequences (683,534 sequences defined as training set in ML Task 1). We then predicted the deleted (missing) data using the multivariate imputation by chained random forests (MICRF) algorithm<sup>79</sup> via the *missRanger* R package<sup>172</sup>. We repeated this step 20 times for each subsample size (50, 100, 500, 1000, 10000, 20000 antibodies) and reported the mean  $R^2$ .

For both ML tasks—and both embedding types implemented in ML Task 1—we performed ablation studies by randomly permuting the column values in the input datasets for the ML models, and confirmed that the prediction accuracy was abolished ( $R^2 \leq 0$ ).

### Antibody sequence encoding with PLM

Antibody sequences were encoded using a Protein Language Model (PLM). PLMs are deep neural networks designed to transform protein sequences into contextual embedding vectors depending on the entirety of the protein sequence. In our experiments, we use the PLM ESM-1v since this model is optimized to predict the effects of mutations on the function of proteins<sup>75</sup>.

To extract a global, fixed-size representation from PLM embeddings, we use a compression scheme based on the Beginning-Of-Sequence (BOS) token, as is customary for LLMs, e.g., for the [CLS] token in ref. 173. BOS tokens are trained to provide a *summarized* representation of the entire protein sequence and are, therefore, a natural choice to represent antibodies.

### Graphical illustrations

We used BioRender.com to create the illustrations in Fig. 1, Fig. 4a, and Fig. 6a. Antibody structural images were produced in PyMOL v2.5.5<sup>174</sup>. We generated the remaining figures in RStudio<sup>142</sup> using the ‘ggplot2’ package<sup>175</sup> and figure panels were aggregated with Adobe Illustrator<sup>176</sup>.

### Statistics and reproducibility

All statistics were calculated using the *rstatix* R package (version 0.7.2). The suitable statistical test used in each analysis is reported accordingly when applicable. The sample size, the number of iterations, replications and data shuffling are also reported for all experiments included. Bias in developability parameter value prediction was avoided by including only non-redundant MWAS DPs in the predictability experiments.

### Experimentally generated native antibody sequences (human IgD, IgK and IgL)

**Human subjects, B-cell isolation and RNA extraction.** Human peripheral blood was obtained from one healthy volunteer. Sample acquisition was approved by the Regional Ethics Committee of South-Eastern Norway (project 6544) and informed consent was obtained. All ethical

regulations relevant to human research participants were followed. Samples were collected in BD Vacutainer® K2 EDTA tubes, and pan B cells were isolated by negative selection using MACSxpress® Whole Blood B Cell Isolation Kit (Miltenyi Biotec). The cells were washed with PBS and remaining erythrocytes were lysed using Red Blood Cell Lysis Solution (Miltenyi Biotec). RNA was extracted using RNeasy Kit (Qiagen). RNA quality and concentration were measured using a Nanodrop spectrophotometer (Thermo Fisher Scientific).

**cDNA synthesis.** 200 ng of RNA was used for cDNA synthesis with 1  $\mu$ l of 100  $\mu$ M isotype-specific reverse primers, 1  $\mu$ l of 10 mM dNTP Mix (Thermo Fisher Scientific) and nuclease-free water up to 14.5  $\mu$ l. Mixture was incubated for 5 min at 65 °C, briefly placed on ice and centrifuged. Subsequently, 4  $\mu$ l of 5X RT buffer (Thermo Fisher Scientific), 0.5  $\mu$ l of RiboLock RNase Inhibitor (Thermo Fisher Scientific), and 1  $\mu$ l Maxima RT (Thermo Fisher Scientific) were added and cDNA synthesis was performed at 50 °C for 30 min with reaction termination at 85 °C for 5 min. The obtained cDNA was purified using MinElute PCR Purification Kit (Qiagen) and eluted in 20  $\mu$ l of EB buffer.

**5' Multiplexing (MPTX) PCR.** 4  $\mu$ l of cDNA were amplified with 0.5  $\mu$ l of 100  $\mu$ M Read2U primer, 1  $\mu$ l of chain-specific 5' forward leader primer mix (Supplementary Data 2), 10  $\mu$ l of KAPA HiFi HotStart ReadyMix (Roche Molecular Systems), and 4.5  $\mu$ l of nuclease-free water at the following conditions: 96 °C—5 min; 25 cycles of 95 °C—20 s, 68 °C—20 s, 72 °C—20 s; 72 °C—5 min; 4 °C—hold. Amplified product was run in a 1.2% agarose gel in TBE buffer. Bands corresponding to the amplified regions of interest (~480 bp) were cut and purified using QIAquick Gel Extraction Kit (Qiagen) with elution in 20  $\mu$ l of EB buffer.

**NGS library generation.** For indexing PCR, 10 ng of the 5' MTPX product were mixed with 0.5  $\mu$ l of 100  $\mu$ M P5\_R1 forward indexing primer, 0.5  $\mu$ l of 100  $\mu$ M P7\_R2 reverse indexing primer containing Illumina index sequence, 12  $\mu$ l of KAPA HiFi HotStart ReadyMix, and nuclease-free water up to 24  $\mu$ l and reaction was performed at the following conditions: 96 °C—5 min; 10 cycles of 95 °C—30 s, 68 °C—30 s, 72 °C—30 s; 72 °C—10 min; 4 °C—hold. The resulting libraries were bead-purified with AMPure XP (Beckman Coulter) using a 1:1 beads ratio. The molarity of the libraries was determined using Qubit™ 4 Fluorometer. The final libraries were inspected with BioAnalyzer (average product length 550 bp) and sequenced on the Illumina MiSeq platform (V3 chemistry 300×2 bp). The raw sequencing data is available on the Sequence Read Archive (BioProject number PRJNA1043047).

### Validation datasets

**The crystal structures dataset.** We obtained 859 crystal structures of paired-chain antibodies (Fv regions) from the antibody structure database (AbDb)<sup>177</sup>. We extracted the amino acid sequences for both (heavy and light) chains and utilized them to further predict the structure of paired and unpaired antibody chains *in silico*. We used several tools to benchmark our analysis including ABodyBuilder<sup>37</sup>, ABodyBuilder2<sup>38</sup>, IgFold<sup>39</sup>, AlphaFold2<sup>178</sup>, and AlphaFold-multimer<sup>179</sup>. Both ABB and ABB2 pipelines include a relaxation step after the initial antibody structure prediction. Of note, the alignment process and regional sequence definition (CDRs and FRs) was performed similarly to the processing steps mentioned for the therapeutic antibodies dataset (see Methods). We finally computed the values of structure-based DPs on these structures to perform the rigid model analysis described in Supplementary Fig. 3.

**The in silico mutated antibody dataset at CDRs.** We subsampled 10 WT antibodies per isotype, except 9 for IgD, and all their corresponding CDR mutants (a total of 30,015 antibody sequences) from the *in silico* mutated dataset (see Methods), and predicted their structures (including WT antibodies) using IgFold version 0.1.0<sup>39</sup>. We used IgFold for this task as

template search-based tools (including ABB) tend to utilize an identical template for antibodies that are within the distance of a single amino acid substitution from their wildtype antibody<sup>180</sup>. We used this dataset to conduct the structural variance study described in Supplementary Fig. 11A.

**The AbDb antibody pairs.** To study the structural variance of antibody mutants with a single aa difference, we sourced 10 pairs of antibodies (Supplementary Table 3) from the public AbDb database<sup>177</sup>. We ensured that each pair of antibodies (1) has the same sequence length and (2) the different (single) amino acid is located in the loops. This dataset was used in the structural variance study described in Supplementary Fig. 11B, C.

### Molecular dynamics simulations

We selected five paired-chain ( $V_H$  and  $V_L$ ) antibody structures from the crystal structures dataset with the best resolutions (Supplementary Table 4) to perform classical MD simulations. Two antibody structures (4TRP and 5WCA) contained missing atoms, which were corrected with the MODELLER's function "complete\_pdb"<sup>181</sup>. Hydrogens were added to initial antibody structures using Reduce<sup>151</sup>. All MD steps were performed in Gromacs v.2022.4<sup>182</sup>. We used AMBER99SB-ILDN<sup>183</sup> force field and transferable intermolecular potential with 3 points<sup>184</sup> water model for MD system preparation<sup>185</sup>. The simulation box was defined as a cube centered around the antibody placed at the 1 nm distance between the solute and the box. All systems contained both Na and Cl ions at 0.1 mol/liter salt concentration. We minimized the energy of each initial system with the steepest descent algorithm<sup>186</sup> for 20000 steps. For the equilibration<sup>187</sup>, we first considered constant volume simulations at 300 K for 100 ps and followed up by constant pressure simulations at 1 bar for another 10 ns. During these equilibration simulations, where applicable, we used the Parrinello-Rahman barostat<sup>188</sup> and V-rescale thermostat<sup>189</sup> with velocity rescaling using 0.1 ps and 0.1 ps time constants, respectively. During all the simulations, we constrained the length of all bonds using the linear constraint solver (LINCS) algorithm<sup>190</sup> and kept the water molecules rigid via the SETTLE algorithm<sup>191</sup>. We used particle mesh Ewald (PME)<sup>192</sup> for treating the electrostatic interactions with a real-space cutoff of 1.0 nm<sup>192</sup>. We simulated a 100 ns independent production run with 2 fs time step for each system, all continued from the last step of the NPT simulations at 300 K and 1 bar using a 2 fs time step. Further details of the parameter specifications can be found in .mdp files. MD simulations were performed starting from crystal structures and ABB-predicted models for the five selected antibodies (Supplementary Table 4), followed by the computation of structure-based DPs on the convergent frames (4001 frames per antibody per structure origin). This data was used to perform the molecular dynamics analysis described in Supplementary Fig. 4.

### Computation of the overlap index ( $\eta$ )

For the analysis in Supplementary Fig. 4C, we used the 'overlapping' R package<sup>193</sup> to quantify the overlap between the distributions of DP values when measured on the MD-simulated crystal structures and the MD-simulated ABB-predicted structures.

### Reporting summary

Further information on research design is available in the Nature Portfolio Reporting Summary linked to this article.

### Data availability

Datasets are available on GitHub: [https://github.com/csi-greifflab/developability\\_profiling](https://github.com/csi-greifflab/developability_profiling). Structural data, such as predicted models and MD trajectories, are stored on Zenodo<sup>194</sup>. The raw sequencing data for the human donor is available on the Sequence Read Archive (BioProject number PRJNA1043047).

### Code availability

Codes are available on GitHub from: [https://github.com/csi-greifflab/developability\\_profiling](https://github.com/csi-greifflab/developability_profiling).

Received: 25 January 2024; Accepted: 5 July 2024;

Published online: 31 July 2024

### References

- Singh, S. et al. Monoclonal antibodies: a review. *Curr. Clin. Pharmacol.* **13**, 85–99 (2018).
- Khetan, R. et al. Current advances in biopharmaceutical informatics: guidelines, impact and challenges in the computational developability assessment of antibody therapeutics. *MABs* **14**, 2020082 (2022).
- Akbar, R. et al. Progress and challenges for the machine learning-based design of fit-for-purpose monoclonal antibodies. *MABs* **14**, 2008790 (2022).
- Laustsen, A. H., Greiff, V., Karatt-Vellatt, A., Muyldermans, S. & Jenkins, T. P. Animal immunization, in vitro display technologies, and machine learning for antibody discovery. *Trends Biotechnol.* **39**, 1263–1273 (2021).
- Wilman, W. et al. Machine-designed biotherapeutics: opportunities, feasibility and advantages of deep learning in computational antibody discovery. *Brief. Bioinform.* **23**, bbac267 (2022).
- Lu, R.-M. et al. Development of therapeutic antibodies for the treatment of diseases. *J. Biomed. Sci.* **27**, 1 (2020).
- Raybould, M. I. J. et al. Five computational developability guidelines for therapeutic antibody profiling. *Proc. Natl Acad. Sci. USA* **116**, 4025–4030 (2019).
- Xu, Y. et al. Structure, heterogeneity and developability assessment of therapeutic antibodies. *MABs* **11**, 239–264 (2019).
- Hummer, A. M., Abanades, B. & Deane, C. M. Advances in computational structure-based antibody design. *Curr. Opin. Struct. Biol.* **74**, 102379 (2022).
- Ahmed, L., Gupta, P. & Martin, K. P. Intrinsic physicochemical profile of marketed antibody-based biotherapeutics. *Proc. Natl Acad. Sci. USA* **118**, e2020577118 (2021).
- Narayanan, H. et al. Design of biopharmaceutical formulations accelerated by machine learning. *Mol. Pharm.* **18**, 3843–3853 (2021).
- Sankar, K. et al. A descriptor set for quantitative structure-property relationship prediction in biologics. *Mol. Inform.* **41**, e2100240 (2022).
- Zarzar, J. et al. High concentration formulation developability approaches and considerations. *MABs* **15**, 2211185 (2023).
- Harmalkar, A. et al. Toward generalizable prediction of antibody thermostability using machine learning on sequence and structure features. *MABs* **15**, 2163584 (2023).
- Zhang, W. et al. Developability assessment at early-stage discovery to enable development of antibody-derived therapeutics. *Antib. Ther.* **6**, 13–29 (2023).
- Carter, P. J. & Lazar, G. A. Next generation antibody drugs: pursuit of the 'high-hanging fruit'. *Nat. Rev. Drug Discov.* **17**, 197–223 (2018).
- Jain, T. et al. Biophysical properties of the clinical-stage antibody landscape. *Proc. Natl Acad. Sci. USA* **114**, 944–949 (2017).
- Evers, A., Malhotra, S. & Sood, V. D. In silico approaches to deliver better antibodies by design: the past, the present and the future. *arXiv* <https://doi.org/10.48550/arXiv.2305.07488> (2023).
- Harvey, E. P. et al. An in silico method to assess antibody fragment polyreactivity. *Nat. Commun.* **13**, 7554 (2022).
- Fernández-Quintero, M. L. et al. Assessing developability early in the discovery process for novel biologics. *MABs* **15**, 2171248 (2023).
- Khan, A. et al. Toward real-world automated antibody design with combinatorial Bayesian optimization. *Cell Rep. Methods* **3**, 100374 (2023).
- Ausserwöger, H. et al. Non-specificity as the sticky problem in therapeutic antibody development. *Nat. Rev. Chem.* **6**, 844–861 (2022).
- Mieczkowski, C. et al. Blueprint for antibody biologics developability. *MABs* **15**, 2185924 (2023).

24. Kingsbury, J. S. et al. A single molecular descriptor to predict solution behavior of therapeutic antibodies. *Sci. Adv.* **6**, eabb0372 (2020).
25. Wolf Pérez, A.-M. et al. In vitro and in silico assessment of the developability of a designed monoclonal antibody library. *MABs* **11**, 388–400 (2019).
26. Han, X., Shih, J., Lin, Y., Chai, Q. & Cramer, S. M. Development of QSAR models for in silico screening of antibody solubility. *MABs* **14**, 2062807 (2022).
27. Widatalla, T., Rollins, Z., Chen, M.-T., Waight, A. & Cheng, A. C. *AbPROP: Language and Graph Deep Learning for Antibody Property Prediction*. [https://icml-compbio.github.io/2023/papers/WCBICML2023\\_paper53.pdf](https://icml-compbio.github.io/2023/papers/WCBICML2023_paper53.pdf) (2023).
28. Licari, G. et al. Embedding dynamics in intrinsic physicochemical profiles of market-stage antibody-based biotherapeutics. *Mol. Pharm.* **2**, 1096–1111 (2022).
29. Makowski, E. K. et al. Optimization of therapeutic antibodies for reduced self-association and non-specific binding via interpretable machine learning. *Nat. Biomed. Eng.* **8**, 45–56 (2023).
30. Jain, T., Boland, T. & Vásquez, M. Identifying developability risks for clinical progression of antibodies using high-throughput in vitro and in silico approaches. *MABs* **15**, 2200540 (2023).
31. Waight, A. B. et al. A machine learning strategy for the identification of key in silico descriptors and prediction models for IgG monoclonal antibody developability properties. *MABs* **15**, 2248671 (2023).
32. Park, E. & Izadi, S. Molecular surface descriptors to predict antibody developability. *bioRxiv* <https://doi.org/10.1101/2023.07.18.549448> (2023).
33. Bauer, J. et al. How can we discover developable antibody-based biotherapeutics? *Front. Mol. Biosci.* **10**, 1221626 (2023).
34. Makowski, E. K. et al. Reduction of monoclonal antibody viscosity using interpretable machine learning. *MABs* **16**, 2303781 (2024).
35. Thrift, W. J. et al. Graph-pMHC: graph neural network approach to MHC class II peptide presentation and antibody immunogenicity. *Brief. Bioinform.* **25**, bbae123 (2024).
36. Raybould, M. I. J., Turnbull, O. M., Suter, A., Guloglu, B. & Deane, C. M. Contextualising the developability risk of antibodies with lambda light chains using enhanced therapeutic antibody profiling. *Commun. Biol.* **7**, 62 (2024).
37. Leem, J., Dunbar, J., Georges, G., Shi, J. & Deane, C. M. ABodyBuilder: Automated antibody structure prediction with data-driven accuracy estimation. *MABs* **8**, 1259–1268 (2016).
38. Abanades, B. et al. ImmuneBuilder: Deep-learning models for predicting the structures of immune proteins. *Commun. Biol.* **6**, 575 (2023).
39. Ruffolo, J. A., Chu, L.-S., Mahajan, S. P. & Gray, J. J. Fast, accurate antibody structure prediction from deep learning on massive set of natural antibodies. *Nat. Commun.* **14**, 2389 (2023).
40. Reynisson, B. et al. Improved prediction of MHC II antigen presentation through integration and motif deconvolution of mass spectrometry MHC eluted. *Ligand Data. J. Proteome Res.* **19**, 2304–2315 (2020).
41. Thorsteinson, N., Gunn, J. R., Kelly, K., Long, W. & Labute, P. Structure-based charge calculations for predicting isoelectric point, viscosity, clearance, and profiling antibody therapeutics. *MABs* **13**, 1981805 (2021).
42. Hutchinson, M. et al. Enhancement of antibody thermostability and affinity by computational design in the absence of antigen. *bioRxiv* <https://doi.org/10.1101/2023.12.19.572421> (2023).
43. Evers, A. et al. Engineering hydrophobicity and manufacturability for optimized biparatopic antibody–drug conjugates targeting c-MET. *MABs* **16**, 2302386 (2024).
44. Sattawa, T. et al. LAP: Liability antibody profiler by sequence & structural mapping of natural and therapeutic antibodies. *PLoS Comput. Biol.* **20**, e1011881 (2024).
45. Feng, J., Jiang, M., Shih, J. & Chai, Q. Antibody apparent solubility prediction from sequence by transfer learning. *iScience* **25**, 105173 (2022).
46. Pudžiuvytė, I. et al. TemStaPro: protein thermostability prediction using sequence representations from protein language models. *Bioinformatics* **40**, btae157 (2024).
47. Manz, R. A., Hauser, A. E., Hiepe, F. & Radbruch, A. Maintenance of serum antibody levels. *Annu. Rev. Immunol.* **23**, 367–386 (2005).
48. Goodnow, C. C., Vinuesa, C. G., Randall, K. L., Mackay, F. & Brink, R. Control systems and decision making for antibody production. *Nat. Immunol.* **11**, 681–688 (2010).
49. Shehata, L. et al. Affinity maturation enhances antibody specificity but compromises conformational stability. *Cell Rep.* **28**, 3300–3308.e4 (2019).
50. Shin, J.-E. et al. Protein design and variant prediction using autoregressive generative models. *Nat. Commun.* **12**, 2403 (2021).
51. Pucca, M. B. et al. History of envenoming therapy and current perspectives. *Front. Immunol.* **10**, 1598 (2019).
52. Krawczyk, K., Raybould, M. I. J., Kovaltuk, A. & Deane, C. M. Looking for therapeutic antibodies in next-generation sequencing repositories. *MABs* **11**, 1197–1205 (2019).
53. Marks, C. & Deane, C. M. How repertoire data are changing antibody science. *J. Biol. Chem.* **295**, 9823–9837 (2020).
54. Petersen, B. M. et al. Regulatory approved monoclonal antibodies contain framework mutations predicted from human antibody repertoires. *Front. Immunol.* **12**, 728694 (2021).
55. Negron, C., Fang, J., McPherson, M. J., Stine, W. B. Jr & McCluskey, A. J. Separating clinical antibodies from repertoire antibodies, a path to in silico developability assessment. *MABs* **14**, 2080628 (2022).
56. Lauer, T. M. et al. Developability index: a rapid in silico tool for the screening of antibody aggregation propensity. *J. Pharm. Sci.* **101**, 2271–2280 (2012).
57. Chen, X. et al. Predicting antibody developability from sequence using machine learning. *bioRxiv* <https://doi.org/10.1101/2020.06.18.159798> (2020).
58. Krawczyk, K., Buchanan, A. & Marcantili, P. Data mining patented antibody sequences. *MABs* **13**, 1892366 (2021).
59. Shetgaonkar, S. & Singh, A. Hybridization of artificial bee colony algorithm with estimation of distribution algorithm for minimum weight dominating set problem. in *ICT Systems and Sustainability* (eds Tuba, M., Akashe, S., Joshi, A.) 607–619 (Springer, Singapore, 2021).
60. Evers, A. et al. SUMO: In silico sequence assessment using multiple optimization parameters. in *Genotype Phenotype Coupling: Methods and Protocols* (eds Zielonka, S. & Krah, S.) 383–398 (Springer US, New York, 2023).
61. Makowski, E. K. et al. Co-optimization of therapeutic antibody affinity and specificity using machine learning models that generalize to novel mutational space. *Nat. Commun.* **13**, 3788 (2022).
62. Saltelli, A. et al. Why so many published sensitivity analyses are false: a systematic review of sensitivity analysis practices. *Environ. Model. Softw.* **114**, 29–39 (2019).
63. Balanda, K. P. & Macgillivray, H. L. Kurtosis: A critical review. *Am. Stat.* **42**, 111–119 (1988).
64. Giudicelli, V. & Lefranc, M. P. Ontology for immunogenetics: the IMGT-ONTOLOGY. *Bioinformatics* **15**, 1047–1054 (1999).
65. Peres, A. et al. IGHV allele similarity clustering improves genotype inference from adaptive immune receptor repertoire sequencing data. *Nucleic Acids Res.* **51**, e86 (2023).
66. Schäfer, J. & Strimmer, K. A shrinkage approach to large-scale covariance matrix estimation and implications for functional genomics. *Stat. Appl. Genet. Mol. Biol.* **4**, Article32 (2005).
67. Barroso, R., Morrison, W. I. & Morrison, L. J. Molecular dissection of the antibody response: opportunities and needs for application in cattle. *Front. Immunol.* **11**, 1175 (2020).

68. Mhanna, V. et al. Adaptive immune receptor repertoire analysis. *Nat. Rev. Methods Prim.* **4**, 1–25 (2024).
69. Sandve, G. K. & Greiff, V. Access to ground truth at unconstrained size makes simulated data as indispensable as experimental data for bioinformatics methods development and benchmarking. *Bioinformatics* **38**, 4994–4996 (2022).
70. Pavlović, M. et al. The immuneML ecosystem for machine learning analysis of adaptive immune receptor repertoires. *Nat. Mach. Intell.* **3**, 936–944 (2021).
71. Perkins, N. J. et al. Principled approaches to missing data in epidemiologic studies. *Am. J. Epidemiol.* **187**, 568–575 (2018).
72. Hong, S. & Lynn, H. S. Accuracy of random-forest-based imputation of missing data in the presence of non-normality, non-linearity, and interaction. *BMC Med. Res. Methodol.* **20**, 199 (2020).
73. Shadbahr, T. et al. The impact of imputation quality on machine learning classifiers for datasets with missing values. *Commun. Med.* **3**, 139 (2023).
74. Alley, E. C., Khimulya, G., Biswas, S., AlQuraishi, M. & Church, G. M. Unified rational protein engineering with sequence-based deep representation learning. *Nat. Methods* **16**, 1315–1322 (2019).
75. Meier, J. et al. Language models enable zero-shot prediction of the effects of mutations on protein function. *bioRxiv* <https://doi.org/10.1101/2021.07.09.450648> (2021).
76. Madani, A. et al. Large language models generate functional protein sequences across diverse families. *Nat. Biotechnol.* **8**, 1099–1106 (2023).
77. Nijkamp, E., Ruffolo, J., Weinstein, E. N., Naik, N. & Madani, A. *ProGen2: Exploring the Boundaries of Protein Language Models*. <https://openreview.net> (2022).
78. Vatcheva, K. P., Lee, M., McCormick, J. B. & Rahbar, M. H. Multicollinearity in regression analyses conducted in epidemiologic studies. *Epidemiology* **6**, 227 (2016).
79. Azur, M. J., Stuart, E. A., Frangakis, C. & Leaf, P. J. Multiple imputation by chained equations: what is it and how does it work? *Int. J. Methods Psychiatr. Res.* **20**, 40–49 (2011).
80. Waljee, A. K. et al. Comparison of imputation methods for missing laboratory data in medicine. *BMJ Open* **3**, e002847 (2013).
81. Aracri, F., Giovanna Bianco, M., Quattrone, A. & Sarica, A. Imputation of missing clinical, cognitive and neuroimaging data of dementia using missForest, a random forest based algorithm. in *2023 IEEE 36th International Symposium on Computer-Based Medical Systems (CBMS)* 684–688 (2023).
82. Molnar, C. *Interpretable Machine Learning*, 318 (Lulu.com, 2020).
83. Teixeira, A. A. R. et al. Simultaneous affinity maturation and developability enhancement using natural liability-free CDRs. *MAbs* **14**, 2115200 (2022).
84. Tiller, T. et al. A fully synthetic human fab antibody library based on fixed VH/VL framework pairings with favorable biophysical properties. *MAbs* **5**, 445–470 (2013).
85. Erasmus, M. F. et al. A single donor is sufficient to produce a highly functional in vitro antibody library. *Commun. Biol.* **4**, 350 (2021).
86. Richardson, E. et al. Characterisation of the immune repertoire of a humanised transgenic mouse through immunophenotyping and high-throughput sequencing. *Elife* **12**, e81629 (2023).
87. Suzek, B. E. et al. UniRef clusters: a comprehensive and scalable alternative for improving sequence similarity searches. *Bioinformatics* **31**, 926–932 (2015).
88. Ofer, D., Brandes, N. & Linial, M. The language of proteins: NLP, machine learning & protein sequences. *Comput. Struct. Biotechnol. J.* **19**, 1750–1758 (2021).
89. Vu, M. H. et al. Linguistically inspired roadmap for building biologically reliable protein language models. *Nat. Mach. Intell.* **5**, 485–496 (2023).
90. Vu, M. H. et al. ImmunoLingo: Linguistics-based formalization of the antibody language. *arXiv* <https://doi.org/10.48550/arXiv.2209.12635> (2022).
91. Schneider, C., Raybould, M. I. J. & Deane, C. M. SAbDab in the age of biotherapeutics: updates including SAbDab-nano, the nanobody structure tracker. *Nucleic Acids Res.* **50**, D1368–D1372 (2022).
92. Bradbury, A. R. M., Dübel, S., Knappik, A. & Plückerthun, A. Animal-versus in vitro-derived antibodies: avoiding the extremes. *MAbs* **13**, 1950265 (2021).
93. Glanville, J. et al. Deep sequencing in library selection projects: what insight does it bring? *Curr. Opin. Struct. Biol.* **33**, 146–160 (2015).
94. Mason, D. M. et al. High-throughput antibody engineering in mammalian cells by CRISPR/Cas9-mediated homology-directed mutagenesis. *Nucleic Acids Res.* **46**, 7436–7449 (2018).
95. Jaffe, D. B. et al. Functional antibodies exhibit light chain coherence. *Nature* **611**, 352–357 (2022).
96. Burbach, S. M. & Briney, B. Improving antibody language models with native pairing. *arXiv* <https://doi.org/10.1016/j.patter.2024.100967> (2023).
97. Akbar, R. et al. A compact vocabulary of paratope-epitope interactions enables predictability of antibody-antigen binding. *Cell Rep.* **34**, 108856 (2021).
98. Norman, R. A. et al. Computational approaches to therapeutic antibody design: established methods and emerging trends. *Brief. Bioinform.* **21**, 1549–1567 (2020).
99. Vishwakarma, P. et al. VHH structural modelling approaches: a critical review. *Int. J. Mol. Sci.* **23**, 3721 (2022).
100. Bailly, M. et al. Predicting antibody developability profiles through early stage discovery screening. *MAbs* **12**, 1743053 (2020).
101. Schoch, A. et al. Charge-mediated influence of the antibody variable domain on FcRn-dependent pharmacokinetics. *Proc. Natl Acad. Sci. USA* **112**, 5997–6002 (2015).
102. Piche-Nicholas, N. M. et al. Changes in complementarity-determining regions significantly alter IgG binding to the neonatal Fc receptor (FcRn) and pharmacokinetics. *MAbs* **10**, 81–94 (2018).
103. Grevys, A. et al. Antibody variable sequences have a pronounced effect on cellular transport and plasma half-life. *iScience* **25**, 103746 (2022).
104. Prihoda, D. et al. BioPhi: A platform for antibody design, humanization, and humanness evaluation based on natural antibody repertoires and deep learning. *MAbs* **14**, 2020203 (2022).
105. Marks, C., Hummer, A. M., Chin, M. & Deane, C. M. Humanization of antibodies using a machine learning approach on large-scale repertoire data. *Bioinformatics* **37**, 4041–4047 (2021).
106. Tennenhouse, A. et al. Computational optimization of antibody humanness and stability by systematic energy-based ranking. *Nat. Biomed. Eng.* **8**, 30–44 (2023).
107. Ramon, A. et al. Assessing antibody and nanobody nativeness for hit selection and humanization with AbNatiV. *Nat. Mach. Intell.* **6**, 74–91 (2024).
108. The Antibody Society. *Antibody Therapeutics Approved or in Regulatory Review in the EU or US*. <https://www.antibodysociety.org/resources/approved-antibodies/> (2022).
109. Tilegenova, C. et al. Dissecting the molecular basis of high viscosity of monospecific and bispecific IgG antibodies. *MAbs* **12**, 1692764 (2020).
110. Seeliger, D. et al. Boosting antibody developability through rational sequence optimization. *MAbs* **7**, 505–515 (2015).
111. Mason, D. M. et al. Optimization of therapeutic antibodies by predicting antigen specificity from antibody sequence via deep learning. *Nat. Biomed. Eng.* **5**, 600–612 (2021).
112. Schretter, C., Kobbelt, L. & Dehaye, P.-O. Golden ratio sequences for low-discrepancy sampling. *J. Graph. Tools* **16**, 95–104 (2012).

113. McKay, M. D., Beckman, R. J. & Conover, W. J. A comparison of three methods for selecting values of input variables in the analysis of output from a computer code. *Technometrics* **21**, 239–245 (1979).
114. Appgar, J. R. et al. Modeling and mitigation of high-concentration antibody viscosity through structure-based computer-aided protein design. *PLoS One* **15**, e0232713 (2020).
115. van der Flier, F. J. et al. What makes the effect of protein mutations difficult to predict? *bioRxiv* <https://doi.org/10.1101/2023.09.25.559319> (2023).
116. Childers, M. C. & Daggett, V. Molecular dynamics methods for antibody design. in *Computer-Aided Antibody Design* (eds Tsumoto, K. & Kuroda, D.) 109–124 (Springer US, 2023).
117. Knapp, B., Frantal, S., Cibena, M., Schreiner, W. & Bauer, P. Is an intuitive convergence definition of molecular dynamics simulations solely based on the root mean square deviation possible? *J. Comput. Biol.* **18**, 997–1005 (2011).
118. Jaszczyszyn, I. et al. Structural modeling of antibody variable regions using deep learning—progress and perspectives on drug discovery. *Front. Mol. Biosci.* <https://doi.org/10.3389/fmolb.2023.1214424> (2023).
119. Kulikova, A. V. et al. Two sequence- and two structure-based ML models have learned different aspects of protein biochemistry. *Sci. Rep.* **13**, 13280 (2023).
120. Makowski, E. K., Chen, H.-T. & Tessier, P. M. Simplifying complex antibody engineering using machine learning. *Cell Syst.* **14**, 667–675 (2023).
121. Fernández-Quintero, M. L. et al. Challenges in antibody structure prediction. *MAbs* **15**, 2175319 (2023).
122. Lane, T. J. Protein structure prediction has reached the single-structure frontier. *Nat. Methods* **20**, 170–173 (2023).
123. Fernández-Quintero, M. L. et al. Characterizing the diversity of the CDR-H3 loop conformational ensembles in relationship to antibody binding properties. *Front. Immunol.* **9**, 3065 (2018).
124. Waibl, F. et al. Comparison of hydrophobicity scales for predicting biophysical properties of antibodies. *Front. Mol. Biosci.* **9**, 960194 (2022).
125. Hoerschinger, V. J. et al. PEP-patch: Electrostatics in protein-protein recognition, specificity, and antibody developability. *J. Chem. Inf. Model.* **63**, 6964–6971 (2023).
126. Eastman, P. et al. OpenMM 7: Rapid development of high performance algorithms for molecular dynamics. *PLoS Comput. Biol.* **13**, e1005659 (2017).
127. Chaudhury, S., Lyskov, S. & Gray, J. J. PyRosetta: a script-based interface for implementing molecular modeling algorithms using Rosetta. *Bioinformatics* **26**, 689–691 (2010).
128. Kang, T. H. & Jung, S. T. Boosting therapeutic potency of antibodies by taming Fc domain functions. *Exp. Mol. Med.* **51**, 1–9 (2019).
129. Abanades, B. et al. The patent and literature antibody database (PLAbDab): an evolving reference set of functionally diverse, literature-annotated antibody sequences and structures. *Nucleic Acids Res.* **52**, D545–D551 (2024).
130. Lee, J., Han, K., Kim, J., Yu, H. & Lee, Y. Solvent: A framework for protein folding. *arXiv* <https://doi.org/10.48550/arXiv.2307.04603> (2023).
131. Singh, R. et al. Learning the language of antibody hypervariability. *bioRxiv* <https://doi.org/10.1101/2023.04.26.538476> (2023).
132. Khade, P. M., Maser, M., Gligorijevic, V. & Watkins, A. M. Mixed structure- and sequence-based approach for protein graph neural networks with application to antibody developability prediction. *bioRxiv* <https://doi.org/10.1101/2023.06.26.546331> (2023).
133. Akbar, R. et al. In silico proof of principle of machine learning-based antibody design at unconstrained scale. *MAbs* **14**, 2031482 (2022).
134. Chen, V. et al. Best practices for interpretable machine learning in computational biology. *bioRxiv* <https://doi.org/10.1101/2022.10.28.513978> (2022).
135. Kovaltsuk, A. et al. Observed antibody space: a resource for data mining next-generation sequencing of antibody repertoires. *J. Immunol.* **201**, 2502–2509 (2018).
136. Vázquez Bernat, N. et al. High-quality library preparation for NGS-based immunoglobulin germline gene inference and repertoire expression analysis. *Front. Immunol.* **10**, 660 (2019).
137. Bolotin, D. A. et al. MiXCR: software for comprehensive adaptive immunity profiling. *Nat. Methods* **12**, 380–381 (2015).
138. Giudicelli, V. et al. From IMGT-ONTOLOGY to IMGT/HighV-QUEST for NGS immunoglobulin (IG) and T cell receptor (TR) repertoires in autoimmune and infectious diseases. *Autoimmun. Infect. Dis.* <https://doi.org/10.16966/2470-1025.103> (2015).
139. Raybould, M. I. J. et al. Thera-SAbDab: the therapeutic structural antibody database. *Nucleic Acids Res.* **48**, D383–D388 (2020).
140. Dunbar, J. & Deane, C. M. ANARCI: antigen receptor numbering and receptor classification. *Bioinformatics* **32**, 298–300 (2016).
141. Osorio, D., Rondon-Villarreal, P. & Torres, R. *Peptides: A Package for Data Mining of Antimicrobial Peptides*. <https://journal.r-project.org/archive/2015/RJ-2015-001/RJ-2015-001.pdf> (2015).
142. RStudio Team. *RStudio: Integrated Development Environment for R*. <http://www.rstudio.com/> (2020).
143. Kelly, S. M., Jess, T. J. & Price, N. C. How to study proteins by circular dichroism. *Biochim. Biophys. Acta* **1751**, 119–139 (2005).
144. Edelhoch, H. Spectroscopic determination of tryptophan and tyrosine in proteins. *Biochemistry* **6**, 1948–1954 (1967).
145. Pace, C. N., Vajdos, F., Fee, L., Grimsley, G. & Gray, T. How to measure and predict the molar absorption coefficient of a protein. *Protein Sci.* **4**, 2411–2423 (1995).
146. Chailyan, A., Marcatili, P. & Tramontano, A. The association of heavy and light chain variable domains in antibodies: implications for antigen specificity: analysis of VH-VL interface in antibodies. *FEBS J.* **278**, 2858–2866 (2011).
147. Eisenberg, D., Weiss, R. M. & Terwilliger, T. C. The hydrophobic moment detects periodicity in protein hydrophobicity. *Proc. Natl Acad. Sci. USA* **81**, 140–144 (1984).
148. Ikai, A. Thermostability and aliphatic index of globular proteins. *J. Biochem.* **88**, 1895–8 (1980).
149. Guruprasad, K., Reddy, B. V. B. & Pandit, M. W. Correlation between stability of a protein and its dipeptide composition: a novel approach for predicting in vivo stability of a protein from its primary sequence. *Protein Eng. Des. Sel.* **4**, 155–161 (1990).
150. Hon, J. et al. SoluProt: prediction of soluble protein expression in *Escherichia coli*. *Bioinformatics* **37**, 23–28 (2021).
151. Word, J. M., Lovell, S. C., Richardson, J. S. & Richardson, D. C. Asparagine and glutamine: using hydrogen atom contacts in the choice of side-chain amide orientation. *J. Mol. Biol.* **285**, 1735–1747 (1999).
152. Brandon, C. J., Martin, B. P., McGee, K. J., Stewart, J. J. P. & Braun-Sand, S. B. An approach to creating a more realistic working model from a protein data bank entry. *J. Mol. Model.* **21**, 3 (2015).
153. Cock, P. J. A. et al. Biopython: freely available python tools for computational molecular biology and bioinformatics. *Bioinformatics* **25**, 1422–1423 (2009).
154. Mitternacht, S. FreeSASA: An open source C library for solvent accessible surface area calculations. *F1000Res.* **5**, 189 (2016).
155. Olsson, M. H. M., Søndergaard, C. R., Rostkowski, M. & Jensen, J. H. PROPKA3: Consistent treatment of internal and surface residues in empirical pKa predictions. *J. Chem. Theory Comput.* **7**, 525–537 (2011).
156. Bakan, A., Meireles, L. M. & Bahar, I. ProDy: protein dynamics inferred from theory and experiments. *Bioinformatics* **27**, 1575–1577 (2011).
157. Pilgrim, M. & Willison, S. *Dive into Python 3* 2nd edn, 412 (Springer, 2009).

158. Jubb, H. C. et al. Arpeggio: A web server for calculating and visualising interatomic interactions in protein structures. *J. Mol. Biol.* **429**, 365–371 (2017).
159. Gu, Z., Eils, R. & Schlesner, M. Complex heatmaps reveal patterns and correlations in multidimensional genomic data. *Bioinformatics* **32**, 2847–2849 (2016).
160. Zhong, J. *Csi-GreiffLab/mwds\_calculator*. [https://github.com/csi-greiffLab/mwds\\_calculator](https://github.com/csi-greiffLab/mwds_calculator) (2023).
161. Shannon, P. et al. Cytoscape: a software environment for integrated models of biomolecular interaction networks. *Genome Res.* **13**, 2498–2504 (2003).
162. Otasek, D., Morris, J. H., Bouças, J., Pico, A. R. & Demchak, B. Cytoscape automation: empowering workflow-based network analysis. *Genome Biol.* **20**, 185 (2019).
163. Kassambara, A. & Mundt, F. *Factoextra: Extract and Visualize the Results of Multivariate Data Analyses*. <https://CRAN.R-project.org/package=factoextra> (2020).
164. R Core Team. *R: A Language and Environment for Statistical Computing*. <https://www.R-project.org/> (2020).
165. Greiff, V. et al. A minimal model of peptide binding predicts ensemble properties of serum antibodies. *BMC Genomics* **13**, 79 (2012).
166. van der Loo, M. P. J. The stringdist package for approximate string matching. *R. J.* **6**, 111–122, <https://CRAN.R-project.org/package=stringdist> (2014).
167. Bachmann, M. *Levenshtein Python Package* <https://pypi.org/project/python-Levenshtein/> (2022).
168. Pedregosa, F. et al. *Scikit-learn: Machine Learning in Python*. <https://scikit-learn.org/stable/> (2011).
169. Dask Development Team. *Dask: Library for dynamic task scheduling*. <https://dask.org> (2016).
170. Edgar, R. C. Search and clustering orders of magnitude faster than BLAST. *Bioinformatics* **26**, 2460–2461 (2010).
171. Kvalseth, T. O. Cautionary note about R<sup>2</sup>. *Am. Stat.* **39**, 279–285 (1985).
172. Mayer, M. *MissRanger: Fast Imputation of Missing Values*. <https://CRAN.R-project.org/package=missRanger> (2023).
173. Devlin, J., Chang, M.-W., Lee, K. & Toutanova, K. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. *arXiv [cs.CL]* (2018).
174. Schrödinger, L. L. C. & DeLano, W. *PyMOL*. <https://pymol.org/> (2024).
175. Wickham, H. Ggplot2. *Wiley Interdiscip. Rev. Comput. Stat.* **3**, 180–185 (2011).
176. Adobe Inc. *Adobe Illustrator*. <https://www.adobe.com> (2019).
177. Ferdous, S. & Martin, A. C. R. AbDb: antibody structure database—a database of PDB-derived antibody structures. *Database* **2018**, bay040 (2018).
178. Jumper, J. et al. Highly accurate protein structure prediction with AlphaFold. *Nature* **596**, 583–589 (2021).
179. Evans, R. et al. Protein complex prediction with AlphaFold-multimer. *bioRxiv* <https://doi.org/10.1101/2021.10.04.463034> (2022).
180. Kuhlman, B. & Bradley, P. Advances in protein structure prediction and design. *Nat. Rev. Mol. Cell Biol.* **20**, 681–697 (2019).
181. Sali, A. & Blundell, T. L. Comparative protein modelling by satisfaction of spatial restraints. *J. Mol. Biol.* **234**, 779–815 (1993).
182. Abraham, M. et al. *GROMACS 2023.1 Manual*. <https://doi.org/10.5281/zenodo.7852189> (2023).
183. Lindorff-Larsen, K. et al. Improved side-chain torsion potentials for the amber ff99SB protein force field. *Proteins* **78**, 1950–1958 (2010).
184. Jorgensen, W. L., Chandrasekhar, J., Madura, J. D., Impey, R. W. & Klein, M. L. Comparison of simple potential functions for simulating liquid water. *J. Chem. Phys.* **79**, 926–935 (1983).
185. Bekker, G.-J., Fukuda, I., Higo, J. & Kamiya, N. Mutual population-shift driven antibody-peptide binding elucidated by molecular dynamics simulations. *Sci. Rep.* **10**, 1406 (2020).
186. Haug, E. J., Arora, J. S. & Matsui, K. A steepest-descent method for optimization of mechanical systems. *J. Optim. Theory Appl.* **19**, 401–424 (1976).
187. Braun, E. et al. Best practices for foundations in molecular simulations. *Living J. Comput. Mol. Sci.* **1**, 5957 (2019).
188. Parrinello, M. & Rahman, A. Crystal structure and pair potentials: a molecular-dynamics study. *Phys. Rev. Lett.* **45**, 1196–1199 (1980).
189. Bussi, G., Donadio, D. & Parrinello, M. Canonical sampling through velocity rescaling. *J. Chem. Phys.* **126**, 014101 (2007).
190. Hess, B., Bekker, H., Berendsen, H. J. C. & Fraaije, J. G. E. M. LINCS: A linear constraint solver for molecular simulations. *J. Comput. Chem.* **18**, 1463–1472 (1997).
191. Miyamoto, S. & Kollman, P. A. Settle: An analytical version of the SHAKE and RATTLE algorithm for rigid water models. *J. Comput. Chem.* **13**, 952–962 (1992).
192. Simmonett, A. C. & Brooks, B. R. A compression strategy for particle mesh Ewald theory. *J. Chem. Phys.* **154**, 054112 (2021).
193. Pastore, M., Loro, P. A. D., Mingione, M. & Calcagni, A. *Overlapping: Estimation of Overlapping in Empirical Distributions*. <https://CRAN.R-project.org/package=overlapping> (2022).
194. Smorodina, E. Structural data for the antibody developability manuscript: cartography of developability landscapes in native and human-engineered antibodies. *Zenodo* <https://doi.org/10.5281/zenodo.10013524> (2023).

## Acknowledgements

Funding was provided by Leona M. and Harry B. Helmsley Charitable Trust (#2019PG-T1D011, to VG), UiO World-Leading Research Community (to VG), UiO: LifeScience Convergence Environment Immunolingo (to VG and GKS), EU Horizon 2020 iReceptor Plus (#825821) (to VG), a Norwegian Cancer Society Grant (#215817, to VG), Research Council of Norway projects (#300740, #311341, #331890 to VG), a Research Council of Norway IKTPLUSS project (#311341, to VG and GKS), and Stiftelsen Kristian Gerhard Jebsen (K.G. Jebsen Coeliac Disease Research Centre, SKGJ-MED-017) (to GKS), and BBSRC (BB/V011065/1, to JG-M). This project has received funding from the Innovative Medicines Initiative 2 Joint Undertaking under grant agreement No 101007799 (Inno4Vac). This Joint Undertaking receives support from the European Union's Horizon 2020 research and innovation programme and EFPIA. This communication reflects the author's view and neither IMI nor the European Union, EFPIA, or any Associated Partners are responsible for any use that may be made of the information contained therein (to VG). We acknowledge the help of Haidara Nadwa (University of Siena) for providing valuable insights in regards to the MD analysis.

## Author contributions

V.G. conceived the study. H.B. and V.G. conceptualized the study and designed the experiments. H.B., E.S., M.P. and J.Z. performed all data analysis and visualization, and developed the project online repositories. H.B., E.S., M.P., J.Z. and V.G. wrote the first draft of the manuscript. H.B. and E.S. prepared sequence and structural developability data (respectively) for all sections of the manuscript. E.S. and R.A. performed antibody structural modeling and molecular dynamic simulations. J.Z. coded and ran the ABC-EDA hybrid algorithm. M.P. and D.N.-Z.G. generated sequence encodings with the protein language model ESM-1v. R.A. and M.C. provided advice on experimental design and statistical analysis. K.L.Q. and I.S. generated the experimental BCR sequences included in this study. K.K. provided the patented antibody dataset (PAD) under a non-disclosure agreement with V.G. P.R., D.N.-Z.G., G.K.S., J.G.-M and J.T.A. contributed through essential scientific discussions and insights for the progress of the project. All authors revised the manuscript and approved its content.

## Competing interests

The authors declare the following competing interests: V.G. declares advisory board positions in aiNET GmbH, Encicom B.V., Absci, Omniscope,



and Diagonal Therapeutics. V.G. is a consultant for Adaptiv Biosystems, Specifica Inc, Roche/Genentech, immunai, Proteinea and LabGenius. H.B. declares a scientific writing post in PipeBio ApS. K.K. is the founder of NaturalAntibody. M.P. and D.N-Z.G. are employed by Adaptiv Biosystems.

### Additional information

**Supplementary information** The online version contains supplementary material available at <https://doi.org/10.1038/s42003-024-06561-3>.

**Correspondence** and requests for materials should be addressed to Habib Bashour or Victor Greiff.

**Peer review information** *Communications Biology* thanks the anonymous reviewers for their contribution to the peer review of this work. Primary Handling Editors: Laura Rodríguez Pérez and Manuel Breuer. A peer review file is available.

**Reprints and permissions information** is available at <http://www.nature.com/reprints>

**Publisher's note** Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

**Open Access** This article is licensed under a Creative Commons Attribution-NonCommercial-NoDerivatives 4.0 International License, which permits any non-commercial use, sharing, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons licence, and indicate if you modified the licensed material. You do not have permission under this licence to share adapted material derived from this article or parts of it. The images or other third party material in this article are included in the article's Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons licence and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this licence, visit <http://creativecommons.org/licenses/by-nc-nd/4.0/>.

© The Author(s) 2024