

# SCIENTIFIC REPORTS



OPEN

## Web-based display of protein surface and pH-dependent properties for assessing the developability of biotherapeutics

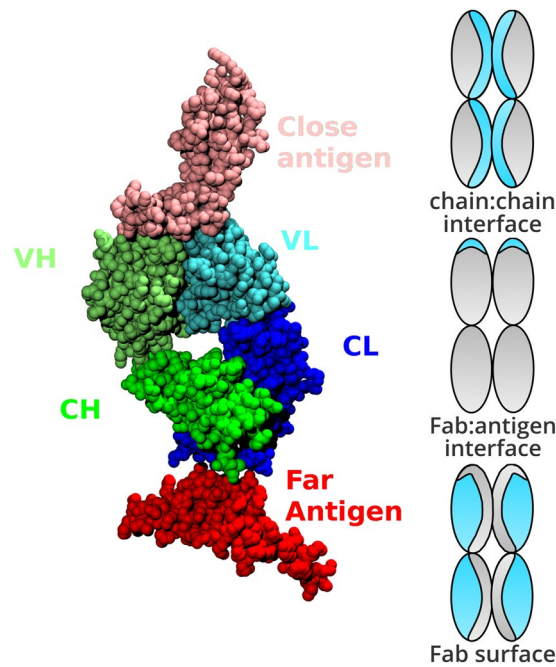
Max Hebditch  & Jim Warwicker 

Protein instability leads to reversible self-association and irreversible aggregation which is a major concern for developing new biopharmaceutical leads. Protein solution behaviour is dictated by the physicochemical properties of the protein and the solution. Optimising protein solutions through experimental screens and targeted protein engineering can be a difficult and time consuming process. Here, we describe development of the protein-sol web server, which was previously restricted to protein solubility prediction from amino acid sequence. Tools are presented for calculating and mapping patches of hydrophobicity and charge on the protein surface. In addition, predictions of folded state stability and net charge are displayed as a heatmap for a range of pH and ionic strength conditions. Tools are evaluated in the context of antibodies, their fragments and interactions. Surprisingly, antibody-antigen interfaces are, on average, at least as polar as Fab surfaces. This benchmarking process provides the user with thresholds with which to assess non-polar surface patches, and possible solubility implications, in proteins of interest. Stability heatmaps compare favourably with experimental data for CH2 and CH3 domains. Display and quantification of surface polarity and pH/ionic strength dependence will be useful generally for investigation of protein biophysics.

Protein biopharmaceuticals (biologics), and in particular monoclonal antibodies, are crucial for many new generation therapeutic interventions<sup>1,2</sup>. Compared to traditional small chemical drugs, antibodies have a higher specificity, as well as target selectivity, leading to fewer off-target effects<sup>3</sup>. However, due to the liquid formulation requirements, and the general instability of proteins compared to small molecules<sup>4,5</sup>, the development of monoclonal antibody biopharmaceuticals can be difficult. Instability of monoclonal antibody products is exacerbated by the delivery requirements. Most biopharmaceutical antibodies are delivered subcutaneously<sup>6</sup>, and this limits the maximum volume to around <1.5 ml, which generally necessitates a concentration of around 100 g/L or higher. This requirement further complicates the delivery of a stable protein formulation, as high concentration often leads to a less stable protein product. Protein instability can lead to non-specific association causing aberrant solution behaviours<sup>7-9</sup>, in more severe cases, instability gives rise to the formation of irreversible, and immunogenic aggregates<sup>10</sup>. Reversible and irreversible association processes limit protein solubility, and have therefore complicated the manufacturing of protein biopharmaceuticals<sup>11-13</sup>. To improve the stability and developability of biologics, various groups have focused on predicting the physicochemical properties of proteins in an attempt to accelerate drug production. Previous work within our group has looked at protein features related to protein solubility, in particular the lack of positively charged surface patches<sup>14</sup>, the ratio of lysine to arginine residues<sup>15</sup>, and the stability of individual Fab domains<sup>16</sup>. Experimental studies<sup>17-25</sup>, as well as computational approaches<sup>26-29</sup>, have been aimed at understanding the solution behaviour of proteins and biologics.

Much research has focused on the role of anisotropic surface patches of charge and hydrophobicity in causing reversible and irreversible protein association<sup>30-35</sup>. This experimental work has led to efforts in predicting protein surface patches *in silico*. For example, the commercially available spatial aggregation propensity (SAP) software<sup>36</sup> which has been applied to IgG antibodies<sup>37</sup>, and is used for predicting aggregation prone hydrophobic regions on the protein surface<sup>38,39</sup>. A development of the SAP software, incorporating charge and hydrophobicity into

School of Chemistry, Manchester Institute of Biotechnology, The University of Manchester, 131 Princess Street, Manchester, M1 7DN, UK. Correspondence and requests for materials should be addressed to J.W. (email: [jim.warwicker@manchester.ac.uk](mailto:jim.warwicker@manchester.ac.uk))



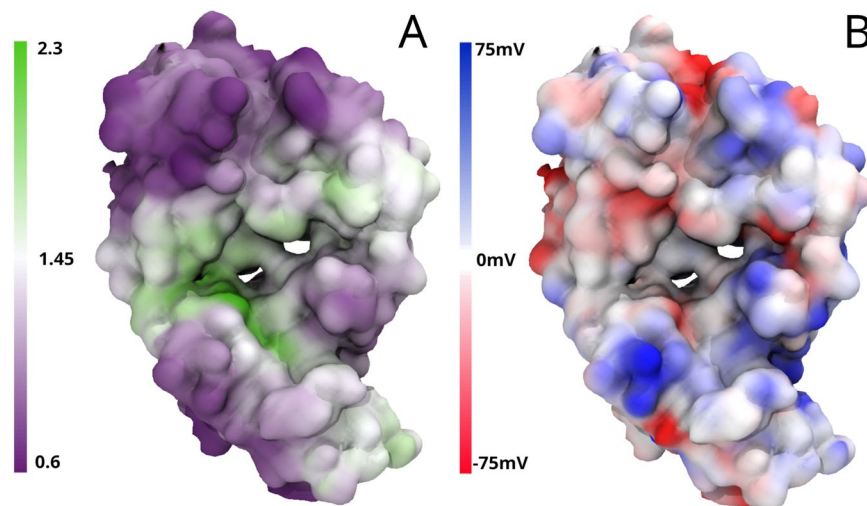
**Figure 1.** Structural classification of the Fab proteins. On the left, an example of categorisation of chains in a Fab PDB file. Although each Fab may contain multiple non Fab chains, we elected to only consider chains within 10Å of the Fab VH and VL domains as the antigen. As a result, the Fab:antigen interface is that between the VH and VL of the Fab, and any non-Fab chain within 10Å. On the right, a schematic representation of classification for atoms in a Fab fragment. Regions highlighted in blue denote interface (top and middle) or surface (non-interface, bottom). Assignment to the different categories was made from calculations of SASA, as described in the text.

the developability index, has been reported<sup>40</sup>. Predicting aggregation risk for antibodies from sequence using bioprocessing data has also been described, with an associated tool available commercially<sup>41</sup>. The freely available CamSol<sup>42,43</sup>, and Aggrescan 3D<sup>44</sup> servers use sequence and structural information for rational design of mutants with enhanced solubility.

We have recently reported<sup>45</sup> the protein-sol server (<https://protein-sol.manchester.ac.uk/>) for sequence-based prediction of protein solubility, calibrated with experimental solubilities in high throughput cell-free expression of *E. coli* proteins. Here, we discuss extension and utility of this freely available web tool, with structure-based calculations. Patch analysis is introduced for electrostatic potential, using Finite Difference Poisson-Boltzmann (FDPB) methods<sup>46</sup> that aid visualisation of asymmetric charge distributions. Analysis of non-polar surface uses a patch approach<sup>47</sup>, importantly with benchmark analysis of Fab fragments to illustrate the range of values that are associated with surfaces and interfaces. Furthermore, taking into account the common use of pH and ionic strength variation in bioprocessing, a heatmap is produced showing prediction of how protein folded state stability varies with these parameters. Comparison with available data for CH2 and CH3 domains reproduces the qualitative differences observed.

## Results

**Categorisation of surface, buried and interface atoms.** Solvent accessible surface area (SASA) was calculated for each atom in each construct: the extracted Fab alone, the extracted antigen alone, each Fab chain alone, and the Fab:antigen complex. With solvent accessible surface areas for each atom in each construct ascertained, we then assigned structural categorisations (Fig. 1) based on solvent accessible surface area. An atom was defined as buried for SASA <5Å<sup>2</sup>, and surface accessible otherwise. A lower threshold (0.1Å<sup>2</sup>) was used to assess change in SASA for an atom, upon interface formation, and assign to the relevant interface (Fab:antigen or Fab chain:Fab chain). Once each atom is tagged with one or two of the above three tags, it was assigned a single structural categorisation, prioritising the interface over surface categorisation. Thus, atoms with a surface categorisation are at the Fab surface and outside of both interface types (Fab:antigen and H chain: L chain). As a result it is possible, for a dataset of Fab fragments, to compare the two interface environments and the remaining surface regions. The ratio of non-polar to polar SASA (NPP ratio) values in an interface are assigned from the constituent parts of the complex that contains that interface. For example, for Fab:antigen, Fab fragment NPP ratio values are taken from the Fab calculation. It is then possible to compare the distributions of NPP ratio for interfacial (including different types of interface) and surface atoms. A similar comparison can be made for distributions reduced to just the set of maximal NPP ratio values, where a maximum is taken from each environment (interface, surface) in each Fab system. For further comparison, the surfaces of a monomeric enzymes set are also included, presumably representative of few interacting surfaces.



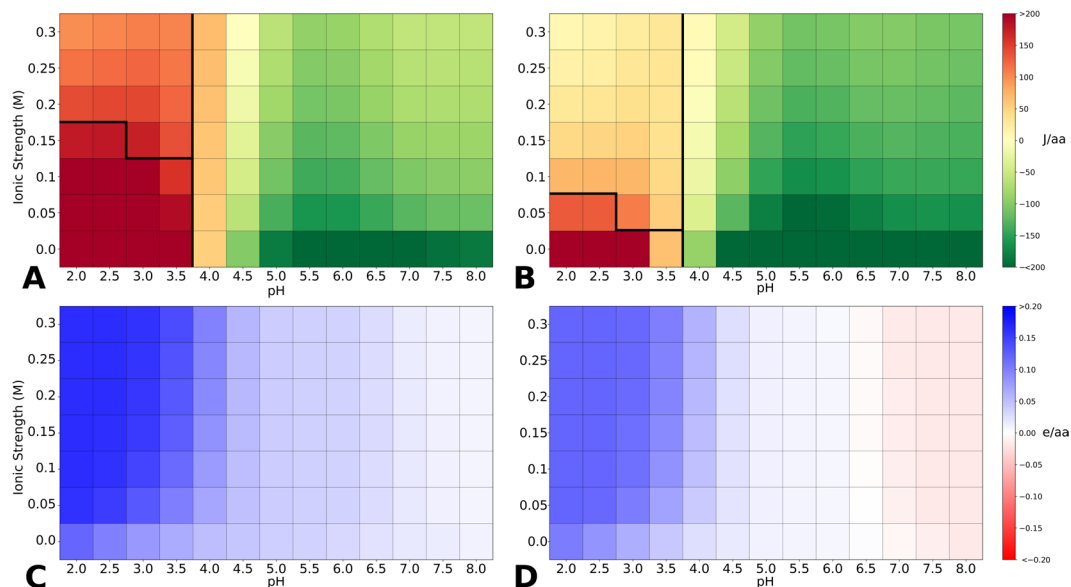
**Figure 2.** Example protein-sol visualisation of surface patches on a Fab. In (A) the Fab is colour-coded from low NPP ratio (purple) to high NPP ratio (green), and in (B) the Fab is colour-coded from negative charge (red) to positive charge (blue). Both are visualised using the embedded NGL viewer on the protein-sol web application after calculation.

**One step visualisation of charged and hydrophobic surface patches.** The protein-sol patches software takes a protein data bank (PDB) structure and calculates patches of charge (at pH 6.3) and hydrophobicity across the protein surface (Fig. 2). This allows the user to quickly identify interesting regions on the protein surface which may influence the behaviour and stability of the protein structure. Electrostatic surface potential based on FDPB calculation is plotted alongside the potential colour-code. Note that a relatively large change in pKa for acidic and basic groups in a surface salt-bridge may be about 1 pH unit, equivalent to 57 mV, referring to the electrostatic potential change due to neighbouring charges that would cause a unit pKa change for a protonating group at 300 K, i.e. the range of potential given here is that expected for electrostatic interactions at the protein surface. Potential equivalent to thermal energy at 300 K,  $kT/e$ , is 25 mV. A scale is also given for colour-coding by patch NPP ratio, from 0.6 (more polar) to 2.5 (more non-polar). Importantly, an additional bar graphic displays the maximum of NPP ratio, in the context of maxima found for interface and non-interface regions of Fab fragments. This information allows the user to find not just the most non-polar region, but also to assess its significance against known interfaces, significantly enhancing the practicality of the tool.

Using our dataset of different Fab structural categorisations, we demonstrate the potential of the protein-sol patches software to quickly identify hotspots of relative hydrophobicity (higher NPP ratio). Through interrogation of the dataset for extremes, Figure S1 shows particularly hydrophobic patches for the Fab chain:chain interface (Figure S1A), a patch on an interface between Fab fragment and antigen (Figure S1B), and a patch on the surface of a Fab fragment (Figure S1C). Note that these are extreme values and, as we show in a subsequent section, antibody-antigen interfaces are, on average, relatively polar.

**Heatmaps show the predicted pH and ionic strength dependence of stability.** Alongside surface visualisation, protein-sol also provides heatmaps for the pH and ionic strength dependence of folded state protein stability, using the Debye-Hückel (DH) method for interactions between ionisable groups, and pKa calculations. Two separate heatmaps (Fig. 3) display predicted charge (units of  $e$  per amino acid), and predicted pH-dependent contribution to stability (J per amino acid). Normalisation relative to sequence length allows direct comparison of proteins. Each heatmap consists of 91 combinations of pH and ionic strength. In order to compare qualitatively with experiment, CH2 and CH3 domains from the IgG1 PDB structure 1HZH are used<sup>48</sup>, since pH and ionic strength variations in stability have been reported for IgG1 CH2 and CH3 domains<sup>49</sup>. Measured phase diagram boundaries for these domains, derived at acidic pH<sup>49</sup> have been marked (Fig. 3), showing a good match between these calculations and experiment. For example, looking across the pH range at 0.15 M ionic strength, there is a greater variation in pH-dependence for CH2 than for CH3 domains. The pH-dependence of stability is directly related to charge<sup>50</sup>, but perhaps the most convenient feature to extract from heatmaps of charge is the predicted sign of net charge. It should be noted that these calculations lack post-translational modifications, such as glycosylation. Here, the IgG1 CH2 domain is predicted to be rather more positively-charged than the CH3 domain at equivalent pH and ionic strength.

Differences in the physicochemical screens represented by the heatmaps can also be visualised in terms of line graph streams that show the variation with pH of charge or energy at a fixed ionic strength (Figure S2). Our dataset of Fab fragments is used as a background with which to compare the pH dependence calculation for the user structure. We chose this background set due to its importance in the biopharmaceutical and biotechnology areas. Wider-scale analysis is possible, but analysis of proteins that are native to differing environments can be a complex topic<sup>51</sup>. Comparison of Figure S2 panels A and B show clearly that the predicted pH-dependence of



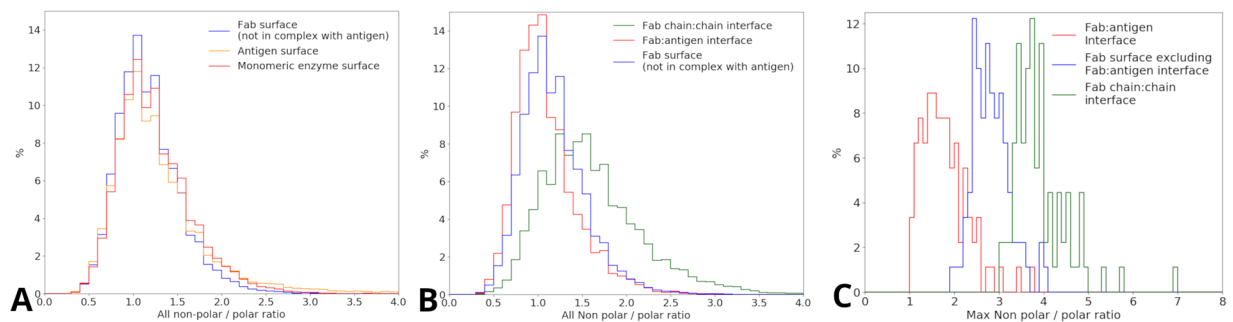
**Figure 3.** Protein-sol stability and charge heatmaps. Calculated folded state stability for the CH2 (A) and CH3 (B) domains, and charge heatmap calculated for the CH2 (C) and CH3 (D) domains of a monoclonal antibody (1HZH).

stability is greater for the CH2 domain for the CH3 domain. Similarly, Figure S2 panels C and D show the more positive predicted charge of the CH2 domain relative to the CH3 domain, at equivalent pH values.

**A common distribution of patch polarity for protein surfaces.** To provide the user with context for patch polarity, alongside surface display, it was necessary to undertake a bioinformatics analysis. Protein surfaces were studied for 3 datasets, the Fab fragments, their corresponding (protein) antigens, and a set of enzymes that are monomeric in their biological states. Surfaces were assigned for Fab fragments and antigens, excluding interfaces, as described in the Methods section, and the entire surfaces of the monomeric enzymes were included. The distribution of NPP ratios are similar for all 3 datasets (Fig. 4A), giving confidence that this form of distribution is broadly representative of protein surfaces. Since a higher NPP ratio reflects a more hydrophobic patch, and a lower NPP ratio relates to a more polar patch, the similarity in distributions suggests that the polar to non-polar spectrum of a protein surface is a general property. All three protein sets have a peak in the distribution at an NPP ratio around 1.0, i.e. with equal contribution from polar and non-polar surface areas.

**The heavy to light chain interface is more non-polar than protein surface, but the Fab-antigen interface is relatively polar.** Having ascertained that the Fab protein surface is representative of protein surfaces in general, a comparison is made with heavy to light chain and Fab-antigen interfaces (Fig. 4B). As expected for a protein-protein interface, the NPP ratio distribution is substantially shifted towards non-polar for the chain-chain interface within a Fab fragment. By comparison, the Fab:antigen interface is slightly more polar than the Fab surface, suggesting that the two interface types have different physicochemical properties. The polarity surface displayed on the server relates to the NPP ratio distributions given in Fig. 4A and B. In order to give the user a more succinct indication of where a protein fits in terms of non-polar surface, we decided to extend the analysis to record the patch with the highest NPP ratio, in a particular protein. It is this property that is shown in the bar chart displayed following the surfaces on the server. Figure 4C shows distributions of the highest NPP ratio values, extracted (one for each protein) from the full distributions (Fig. 4B). Interfaces within a Fab fragment are again, on average, shifted towards more non-polar. The relative polarity of Fab-antigen interfaces is now even more apparent, clearly shifted towards more polar than Fab surface. For each Fab, the most hydrophobic patch is generally in the the Fab chain:chain interface, and the least hydrophobic in the Fab:antigen interface. Overall, this approach allows a user to find the single largest non-polar patch within a structure, and is typically related to features used in the assessment of the developability for biotherapeutics.

Statistical comparisons of the distributions in Fig. 4 are given in Table S1. The means for the (all) surface distributions of NPP ratio in Fig. 4A are close (Fab surface 1.28, Enzyme surface 1.35, Antigen surface 1.22). Significant differences (Table S1) between the distributions displayed in Fig. 4A originate from variation in the tails at higher NPP ratio. In Fig. 4 panels B and C, significant differences (Table S1) underpin the separation of interface and surface distributions for both all and maximal NPP ratios. In wider comparison to other methods of displaying and quantifying surface properties of proteins, it should be emphasised that proteins are typically displayed with the solvent re-entrant surface, as is the case in our server. However, the property we are displaying on that surface is a patch average of NPP ratio over the solvent accessible surface, since the patch size approximates to a typical protein-protein interface area, and hydrophobic interactions are thought to be mediated by extent of non-polar SASA rather than molecular surface.



**Figure 4.** Distribution of hydrophobic patches across the different Fab structural classifications. **(A)** Distributions of all NPP ratio values for the Fab, antigen and monomeric enzyme surfaces combined for each protein. **(B)** Distribution of all NPP values for all three Fab structural categorisations: Fab:antigen interface, Fab surface (excluding Fab:antigen interface) and Fab chain:chain interface combined for each protein and compared. **(C)** Distribution of max NPP ratio values for Fab:antigen interface, Fab chain:chain interface, and Fab surface.

## Discussion

The development and use of therapeutic proteins can be limited by instabilities which complicate manufacture, storage and delivery. It is important to improve understanding and to provide predictions for the factors that cause reversible and irreversible association. To help improve the developability of biopharmaceuticals, in past work, we introduced the protein-sol sequence software for predicting protein solubility based on primary structure<sup>45</sup>. In this work, we introduce two new tools, available for free and with no licensing requirements, protein-sol patches and heatmaps. Whilst targeted at the biopharmaceutical research community, they could also be of wider interest for biotechnology. Both hydrophobic and charged surface patches have been implicated in aberrant solution behavior. The protein-sol patches code is used to calculate the predicted surface patches from a protein structure. Incorporation of the NGL viewer allows fast and simple visualisation of the surface electrostatic potential and polarity (hydrophobicity). Calculation results are also available for download, as PDB format files with surface patches colour-coded using the B-factor field. Results are therefore readily available for further processing or visualisation by the user.

To put the server output into context, we investigated the hydrophobic properties of surface and interfacial regions of Fab fragments, as well as a dataset of soluble monomeric enzymes. There is little difference in the means of NPP ratio distributions for surface regions of Fab fragments, their corresponding protein antigens, and monomeric enzymes (Fig. 4A). One interpretation of this result is that a standard protein surface profile of polarity is associated with a balance between structural stability and solubility. As expected, the heavy - light chain interface is relatively non-polar, both in overall distribution (Fig. 4B) and as peak values for each Fab (Fig. 4C). Surprisingly though, Fab-antigen interfaces are relatively polar, an observation that is exaggerated when viewed as peak NPP ratios (Fig. 4C), as compared with the whole distributions (Fig. 4B). Antibody-antigen interfaces have been reported to differ from other protein-protein interfaces, tending to be smaller in size, incorporating fewer helices and more loops, with less hydrophobic packing<sup>52</sup>. We now find that antibody-antigen interfaces are, if anything, even more polar than surface (non-interfacial) regions. It is possible that the constraints of altered interfacial size and different secondary structure composition lead to a reduction in non-polarity. An alternative possibility is that as a non-obligate interface, there remains a requirement for solubility in the absence of interface formation. Further, following reports that proteins at high naturally occurring concentrations tend to more soluble<sup>16,53</sup>, the constraint for relatively polar surface would be enhanced for antibodies. The observation of relatively polar antibody-antigen interfaces is not necessarily inconsistent with the finding that  $\pi$ - $\pi$  interactions are common<sup>52</sup>. Differences were computed between amino acid percentage compositions overall in Fab fragments and specifically in CDRs, for the dataset used in this study. The largest such difference is for tyrosine, at 5.9%, whilst the other aromatic sidechain amino acid differences are much lower, at 1.5% (phenylalanine) and 0.4% (tryptophan). Next largest after tyrosine is valine (-4.4%). Thus some compensatory effect may exist, with elevation of tyrosine (and the potential for  $\pi$ - $\pi$  interactions) countered by lowering of hydrophobic valine sidechains.

The NPP ratio part of the protein-sol server has been developed to allow users to view non-polar patches in the context of developability. Incorporation of modal values from the NPP peak values distributions in Fig. 4C, alongside the NPP peak value for the user's protein in a simple graphic, will aid such assessment. Whilst non-polar patches may be the focus for developability, we do not discard more polar regions from the display. A demonstration of the relative polarity of antibodies at the antigen-combining site is shown in Figure S3. Whereas for Figure S1B we chose the most non-polar complementarity-determining region (CDR) (with a green patch), more typically the relative polarity is apparent (Figure S3). Indeed, since this graphical tool makes the separation of relatively polar and non-polar regions readily apparent, it could accelerate discoveries such as that made here for antibody - antigen interfaces. The composition of F,W,Y amino acids in the CDRs of the Fab shown in Figure S3 (3mly) is 14.8%, within a range of 10% to 29% for CDRs in the set of Fabs used. By contrast, the F,W,Y composition for the entire Fabs is 8.7% for 3mly, and range from 7% to 12% for the set. Our analysis of surface non-polarity with a 13Å radius patch does not, in general, highlight CDR surfaces. Figure S4 shows the effect of reducing patch radius for 3mly, non-polar patches do now appear within the CDRs, but these remain smaller non-polar features than seen elsewhere on the Fab surface. We infer that the use in CDRs of amino acids with



aromatic sidechains is more subtle than can be judged simply by scale of exposed non-polar surface area. It is possible that non-polar SASA buried at an interface is a poor indicator of hydrophobic contribution to binding energy for antibody-antigen complexation.

The current work provides an additional tool for groups looking to identify regions of a protein for engineering improved solution properties. Both hydrophobic patches<sup>54</sup>, and charged patches<sup>30–32</sup> have been mutated to alter solution behavior. The ability to rapidly visualise surface patches will further inform and accelerate such work.

With regard to heatmap depiction of predicted stability for pH and ionic strength variation, these are two important factors when studying proteins. Although the trend of changes will be uniform, acidification tends towards a more positive protein and increased ionic strength reduces electrostatic interactions, the net outcome is a delicate balance of the constituent parts. For example, we demonstrate (Fig. 3) that qualitative experimentally-determined differences between IgG1 CH2 and CH3 domains<sup>49</sup>, are reproduced by our calculations. Furthermore, the user can view the size of predicted pH-dependence as a comparison with the dataset of Fab fragments, with plots normalized for sequence length. Predicted protein charge is also presented, in an analogous manner. Fab fragments were used as the control dataset since they are a widely used biopharmaceutical platform. We suggest that the protein-sol heatmap may be a useful tool for accelerating formulation screens by identifying potentially favourable regions prior to formulation development, when a structure or structural model is available. Since viral clearance procedures often involve a low pH step, the heatmap analysis will aid determination of the degree to which protein stability is diminished as salt-bridges and other favourable interactions are lost at acidic pH.

In this work we have discussed how the polar, non-polar, pH and ionic strength dependent properties influence protein stability in solution, and how instability can limit development. The protein-sol software suite, incorporating patches and heatmap software has been designed to benefit researchers interested in understanding the surface properties, and stability, of proteins in solution. While developing the server, we have demonstrated how the patches software can identify interesting physicochemical properties of Fab chain:chain and Fab:antigen interfaces, and also how predictions for protein stability compare favourably to measured data. Protein-sol is freely available. Our initial work suggests that it could contribute to the acceleration of protein engineering and formulation optimisation, and to the improvement of developability for new biotherapeutic leads. It also provides insight into the fundamental properties of proteins in solution.

## Methods

**Using the protein-sol patches and heatmap tools.** All software at protein-sol is free to use without license or registration and is available online at <https://protein-sol.manchester.ac.uk>. To use the protein-sol patches or heatmap code, the user simply needs to upload a protein structure in the standard protein data bank format<sup>55</sup>, with results returned in a molecular graphics viewer and as downloadable files, available using a supplied custom URL for 7 days. The protein-sol webserver is built using open source software. Patches data is displayed using the NGL viewer<sup>56</sup>, and the heatmap visualisations are made in python.

**Protein-sol patches calculation and visualisation.** From the supplied PDB structure, only protein is included in the calculation, with the advantage that parameterisation failures for ligands unknown to the dictionary are avoided. Electrostatic calculations follow published protocols from our group<sup>57,58</sup> (Table S2 and S3), but pKa calculations are not made. Ionisable group charges are fixed at pH 6.3, giving half protonation for histidine sidechains, full deprotonation of aspartate and glutamate sidechains and protein carboxy termini, and full protonation of lysine and arginine sidechains and protein amino termini. Any supplied hydrogen atoms are removed, and polar hydrogen atoms added back, to carry partial charges. Electrostatic potential is calculated with the FDPB method<sup>57</sup> on a 0.6 Å spaced grid, with relative dielectric values of 4 for protein and 78.4 for water. Counterions are included at 0.15 M concentration to model ionic strength that matches physiological. A boundary condition of zero potential is set at 10 Å beyond the protein, with the surface between protein and solvent defined by the solvent re-entrant surface. This 10 Å boundary condition is sufficient since at 0.15 M ionic strength, the Debye-Hückel screening length is 8 Å, and potential is displayed at the protein molecular surface rather than being contoured in solvent. For ease of display, both in the server and as a download, the potential is transcribed from the Cartesian calculation grid to the B-factor field of a PDB file containing the original coordinates. To accomplish this with visualisation of potential values at the protein surface, a grid shell surrounding the protein is extracted from the Cartesian grid<sup>59</sup>, and potential values assigned to protein atoms according to the closest point on this surface grid shell. Potential values are capped at lower and upper values of −86 and +86 mV, to fit with the PDB B-factor field, these values correspond to an interaction magnitude for a unit charge in the field of about 8.6 kJ/mole. The resulting electrostatic potential surface, and patches, can be manipulated by the user with the NGL viewer<sup>56</sup>. An equivalent colour scheme for potential can also be viewed from the downloadable coordinate file, for example using the red\_white\_blue spectrum command, with minima and maxima of −86 and +86 in PyMOL<sup>60</sup>. In the embedded viewer of the server, various representations other than surface are possible, as are full-screen viewing and picture download.

A different branch of the code evaluates the non-polarity of patches around over the protein surface. For this purpose, a patch is associated with each non-hydrogen atom in the protein. Each patch is the ratio of non-polar to polar solvent accessible surface area for all non-hydrogen atoms within a 13 Å radius of the central atom<sup>47</sup>. As for electrostatic potential, this property is inserted into the B-factor field of a PDB file, and displayed in the embedded NGL viewer, as well as being downloadable for local viewing. Colour-coding in the embedded viewer is chosen as purple (more polar) to green (more non-polar), so as to distinguish it from the standard red, white, blue scheme for electrostatic potential, and can be visualised using the magenta\_white\_green spectrum command with minima and maxima of 0.4 and 2.5 in PyMOL<sup>60</sup>. It is standard practice to display the most non-polar surface regions,

in the context of protein solubility. It is worth noting however that more polar regions also carry information, to our surprise we found that the antigen-combining regions of antibodies are, on average, relatively polar. A larger SASA threshold ( $5\text{\AA}^2$ ) is used to determine surface or buried atoms, than for identifying atoms at an interface (change of  $0.1\text{\AA}^2$  upon complexation). This difference ensures that atoms that are almost entirely buried are not erroneously labelled as surface. All carbon atoms are assigned to non-polar, and oxygen, nitrogen and sulphur to polar. In practice, sulphur atoms, which could in principle also be assigned as non-polar (oxidised in disulphide bonds), are relatively scarce.

**Protein-sol heatmaps for the predicted pH and ionic strength dependence of stability.** Whereas the calculation of electrostatic potential on protein-sol is made with standard pKas, it is the differences to standard pKas ( $\Delta$  pKas) that determines the pH-dependent contribution to folded state stability. Not only are pH and ionic strength screens used in formulation studies, but also low pH is used for viral inactivation of biologics expressed from CHO cells<sup>61</sup>. We have developed software in previous applications to compute pKas and the pH-dependent contribution to protein stability, and now provide this code at the protein-sol site. Experimental groups often provide the results of pH and ionic strength screens as heatmaps, and we have therefore chosen this format. Whilst we are unlikely to be describing precisely a feature measured experimentally, folded state stability (for which we provide a prediction of the pH-dependent component) is a key underlying property. Indeed, in the Results section we discuss a qualitative fit between heatmaps generated for CH2 and CH3 domains, and experimental data. Rather than the FDPB model used for electrostatic surface generation, we use the more simple Debye-Hückel (DH) scheme for charge-charge interactions in a medium of uniform relative dielectric (78.4, water) and ionic strength (variable, 0 to 0.3 M). The DH calculations of pKas are suitable for systems where the focus is on solvent exposed ionisable groups, such that water solvation and counter-ions act to reduce interactions between ionisable groups<sup>62</sup>. Areas where DH calculations of pH-dependence are likely to be less reliable include enzyme active sites and ion channels, where water is excluded to a large degree. Another approximation in these calculations is the use of a model counterion screening, with no account taken of differences between salt species. Modelling of such differences will require more detailed analysis of specific binding between proteins and counterions. This allows rapid calculation<sup>62</sup> of the required Monte Carlo sampling of protonation states<sup>63</sup>. It is approximated that there are no interactions between ionisable groups in the unfolded state, and the pH-dependent energy is given in Joules per amino acid, a normalisation against protein size. The predicted net charge of the protein (units of e per amino acid) is also given in the heatmap format, within the pH 2 to 8, and ionic strength 0 to 0.3 M ranges. Ligands are, again, excluded from the calculations. In order to give the user context, 2D plots of pH-dependent contribution to stability are drawn for ionic strengths of 0, 0.15, and 0.3 M. The user-supplied protein is displayed against a background of the Fabs dataset analysed in this work.

For calculation with representative CH2 and CH3 domains of an IgG1 antibody, the 1HZH structure was used<sup>48</sup>. Coordinate files for each domain (CH2 and CH3) were extracted from the overall 1HZH file.

**Datasets of Fab structures for calculation.** The Fab dataset was formed by searching the PDB<sup>55</sup> for structures containing Fabs, and the biological assembly files retrieved. Sequences from the resulting structures were analysed manually to identify only structures with unique heavy and light chains, resulting in 199 Fab structures. From these Fab structures, the four individual domains, the variable and constant domains of the heavy chain (VH and CH), and the variable and constant domains of the light chain (VL and CL), were identified using interdomain sequence motifs<sup>45</sup>. This resulted in 199 Fab structures that constitute the heatmap dataset.

In order to identify antibody:antigen interfaces, coordinates for the antigen binding VH and VL domains were compared with all non-Fab atoms in the relevant PDB file. For any non-Fab coordinate within  $10\text{\AA}$  of the combined VH and VL domain coordinates, the entire chain of the close non-Fab structure was extracted and combined with the entire Fab, in a new coordinate file. From the original 199 entirely unique Fab structures, 90 of the original biological assemblies contained an antigen within  $10\text{\AA}$  of the VH and VL domains (Fig. 1). These 90 Fabs were then also split into heavy and light chains for chain:chain analysis, and forming the basis for putting the patches part of the server into the context of Fab calculations. We were interested in whether regions of Fab that were not determined to be interfacial (H - L chain or with antigen) were representative more generally of protein surfaces. For this purpose we used a dataset of 54 enzymes<sup>59</sup> known to be monomeric, and thus likely to present mostly non-interfacial amino acids. Statistical comparison of NPP ratio distributions was calculated using the independent two-group Mann-Whitney U test in the statistical programming language R.

## Data Availability

The datasets analysed during this study are available from the corresponding author on reasonable request. The reported web tool is freely available online.

## References

1. Carter, P. J. Introduction to current and future protein therapeutics: A protein engineering perspective. *Exp. Cell. Res.* **317**, 1261–1269 (2011).
2. Ecker, D. M., Jones, S. D. & Levine, H. L. The therapeutic monoclonal antibody market. *mAbs* **7**, 9–14 (2015).
3. Smith, A. J. New Horizons in Therapeutic Antibody Discovery: Opportunities and Challenges versus Small-Molecule Therapeutics. *J. Biomol. Screen.* **20**, 437–453 (2014).
4. Wang, W. Instability, stabilization, and formulation of liquid protein pharmaceuticals. *Int. J. Pharm.* **185**, 129–188 (1999).
5. Manning, M. C., Chou, D. K., Murphy, B. M., Payne, R. W. & Katayama, D. S. Stability of protein pharmaceuticals: an update. *Pharm. Res.* **27**, 544–575 (2010).
6. Narasimhan, C., Mach, H. & Shameem, M. High-dose monoclonal antibodies via the subcutaneous route: challenges and technical solutions, an industry perspective. *Ther. Deliv.* **3**, 889–900 (2012).
7. Woods, J. M. & Nesta, D. Formulation effects on opalescence of a high-concentration MAb. *Bioprocess Int.* **8**, 48–59 (2010).

8. Liu, J., Nguyen, M. D. H., Andya, J. D. & Shire, S. J. Reversible self-association increases the viscosity of a concentrated monoclonal antibody in aqueous solution. *J. Pharm. Sci.* **94**, 1928–1940 (2005).
9. Raut, A. S. & Kalonia, D. S. Pharmaceutical Perspective on Opalescence and Liquid-Liquid Phase Separation in Protein Solutions. *Mol. Pharm.* **13**, 1431–1444 (2016).
10. Hansel, T. T., Kropshofer, H., Singer, T., Mitchell, J. A. & George, A. J. T. The safety and side effects of monoclonal antibodies. *Nat. Rev. Drug Discov.* **9**, 325–338 (2010).
11. Shire, S. J. Formulation and manufacturability of biologics. *Curr. Opin. Biotechnol.* **20**, 708–714 (2009).
12. Daugherty, A. L. & Mersny, R. J. Formulation and delivery issues for monoclonal antibody therapeutics. *Adv. Drug Deliv. Rev.* **58**, 686–706 (2006).
13. Mitragotri, S., Burke, P. A. & Langer, R. Overcoming the challenges in administering biopharmaceuticals: formulation and delivery strategies. *Nat. Rev. Drug Discov.* **13**, 655–672 (2014).
14. Chan, P., Curtis, R. & Warwicker, J. Soluble expression of proteins correlates with a lack of positively-charged surface. *Sci. Rep.* **3**, 3333 (2013).
15. Warwicker, J., Charonis, S. & Curtis, R. Lysine and arginine content of proteins: Computational analysis suggests a new tool for solubility design. *Mol. Pharm.* **11**, 294–303 (2014).
16. Hebditch, M., Curtis, R. & Warwicker, J. Sequence composition predicts immunoglobulin superfamily members that could share the intrinsically disordered properties of antibody CH1 domains. *Sci. Rep.* **7**, 12404 (2017).
17. Chari, R., Jerath, K., Badkar, A. V. & Kalonia, D. S. Long- and Short-Range Electrostatic Interactions Affect the Rheology of Highly Concentrated Antibody Solutions. *Pharm. Res.* **26**, 2607–2618 (2009).
18. Esfandiary, R., Parupudi, A., Casas-Finet, J., Gadre, D. & Sathish, H. Mechanism of Reversible Self-Association of a Monoclonal Antibody: Role of Electrostatic and Hydrophobic Interactions. *J. Pharm. Sci.* **104**, 577–586 (2015).
19. Neergaard, M. S. *et al.* Viscosity of high concentration protein formulations of monoclonal antibodies of the IgG1 and IgG4 subclass - Prediction of viscosity through protein-protein interaction measurements. *Eur. J. Pharm. Sci.* **49**, 400–410 (2013).
20. Yearley, E. J. *et al.* Small-angle neutron scattering characterization of monoclonal antibody conformations and interactions at high concentrations. *Biophys. J.* **105**, 720–731 (2013).
21. Calero-Rubio, C., Ghosh, R., Saluja, A. & Roberts, C. J. Predicting protein-protein interactions of concentrated antibody solutions using dilute solution data and coarse-grained molecular models. *J. Pharm. Sci.* **107**, 1269–1281 (2017).
22. Roberts, D. *et al.* Specific ion and buffer effects on protein-protein interactions of a monoclonal antibody. *Mol. Pharm.* **12**, 179–193 (2014).
23. Ghosh, R., Calero-Rubio, C., Saluja, A. & Roberts, C. J. Relating Protein-Protein Interactions and Aggregation Rates from Low to High Concentrations. *J. Pharm. Sci.* **105**, 1086–1096 (2016).
24. Inouye, H., Houde, D., Temel, D. B. & Makowski, L. Utility of Solution X-Ray Scattering for the Development of Antibody Biopharmaceuticals. *J. Pharm. Sci.* **105**, 3278–3289 (2016).
25. Schermeyer, M. T., Wöll, A. K., Kokke, B., Eppink, M. & Hubbuch, J. Characterization of highly concentrated antibody solution - A toolbox for the description of protein long-term solution stability. *mAbs* **9**, 1169–1185 (2017).
26. Calero-Rubio, C., Saluja, A. & Roberts, C. J. Coarse-Grained Antibody Models for "weak" Protein-Protein Interactions from Low to High Concentrations. *J. Phys. Chem. B* **120**, 6592–6605 (2016).
27. Lilyestrom, W., Yadav, S., Shire, S. J. & Scherer, T. M. Monoclonal antibody self-association, cluster formation, and rheology at high concentrations. *J. Phys. Chem. B* **117**, 6373–6384 (2013).
28. Corbett, D. *et al.* Coarse-grained modeling of antibodies from small-angle scattering profiles. *J. Phys. Chem. B* **121**, 8276–8290 (2017).
29. Kuhn, A. B. *et al.* Improved Solution State Properties of Monoclonal Antibodies by Targeted Mutations. *J. Phys. Chem. B* **121**, 10818–10827 (2017).
30. Yadav, S., Shire, S. J. & Kalonia, D. S. Viscosity behavior of high-concentration monoclonal antibody solutions: Correlation with interaction parameter and electroviscous effects. *J. Pharm. Sci.* **101**, 998–1011 (2012).
31. Perchiacca, J. M., Ladiwala, A. R. A., Bhattacharya, M. & Tessier, P. M. Aggregation-resistant domain antibodies engineered with charged mutations near the edges of the complementarity-determining regions. *Protein Eng. Des. Sel.* **25**, 591–601 (2012).
32. Chow, C. K., Allan, B. W., Chai, Q., Atwell, S. & Lu, J. Therapeutic Antibody Engineering to Improve Viscosity and Phase Separation Guided by Crystal Structure. *Mol. Pharm.* **13**, 915–923 (2016).
33. Li, W., Persson, B. A., Lund, M., Bergenholtz, J. & Zackrisson-Oskolkova, M. Concentration-Induced Association in a Protein System Caused by a Highly Directional Patch Attraction. *J. Phys. Chem. B* **120**, 8953–8959 (2016).
34. Roberts, D. *et al.* The role of electrostatics in protein-protein interactions of a monoclonal antibody. *Mol. Pharm.* **11**, 2475–2489 (2014).
35. Austerberry, J. I. *et al.* The effect of charge mutations on the stability and aggregation of a human single chain Fv fragment. *Eur. J. Pharm. Biopharm.* **115**, 18–30 (2017).
36. Chennamsetty, N., Voynov, V., Kayser, V., Helk, B. & Trout, B. L. Design of therapeutic proteins with enhanced stability. *Proc. Natl. Acad. Sci.* **106**, 11937–11942 (2009).
37. Chennamsetty, N., Helk, B., Voynov, V., Kayser, V. & Trout, B. L. Aggregation-prone motifs in human immunoglobulin G. *J. Mol. Biol.* **391**, 404–413 (2009).
38. Courtois, F., Agrawal, N. J., Lauer, T. M. & Trout, B. L. Rational design of therapeutic mAbs against aggregation through protein engineering and incorporation of glycosylation motifs applied to bevacizumab. *mAbs* **8**, 99–112 (2016).
39. Voynov, V., Chennamsetty, N., Kayser, V., Helk, B. & Trout, B. L. Predictive tools for stabilization of therapeutic proteins. *mAbs* **1**, 580–582 (2009).
40. Lauer, T. M. *et al.* Developability index: a rapid in silico tool for the screening of antibody aggregation propensity. *J. Pharm. Sci.* **101**, 102–115 (2012).
41. Obrezanova, O. *et al.* Aggregation risk prediction for antibodies and its application to biotherapeutic development. *mAbs* **7**, 352–363 (2015).
42. Sormanni, P., Aprile, F. A. & Vendruscolo, M. The CamSol method of rational design of protein mutants with enhanced solubility. *J. Mol. Biol.* **427**, 478–490 (2015).
43. Sormanni, P., Amery, L., Ekizoglou, S., Vendruscolo, M. & Popovic, B. Rapid and accurate in silico solubility screening of a monoclonal antibody library. *Sci. Rep.* **7**, 8200 (2017).
44. Zambrano, R. *et al.* Aggrescan3d (a3d): server for prediction of aggregation properties of protein structures. *Nucleic acids research* **43**, W306–W313 (2015).
45. Hebditch, M., Carballo-Amador, M. A., Charonis, S., Curtis, R. & Warwicker, J. Protein-Sol: A web tool for predicting protein solubility from sequence. *Bioinformatics* **33**, 3098–3100 (2017).
46. Warwicker, J. Continuum dielectric modelling of the protein-solvent system, and calculation of the long-range electrostatic field of the enzyme phosphoglycerate mutase. *J. Theor. Biol.* **121**, 199–210 (1986).
47. Cole, C. & Warwicker, J. Side-chain conformational entropy at protein-protein interfaces. *Protein Sci.* **11**, 2860–2870 (2002).
48. Saphire, E. O. *et al.* Crystal Structure of a Neutralizing Human IgG Against HIV-1: A Template for Vaccine Design. *Science* **293**, 1155–1159 (2001).



49. Yageta, S., Lauer, T. M., Trout, B. L. & Honda, S. Conformational and Colloidal Stabilities of Isolated Constant Domains of Human Immunoglobulin G and Their Impact on Antibody Aggregation under Acidic Conditions. *Mol. Pharm.* **12**, 1443–1455 (2015).
50. Antosiewicz, J., McCammon, J. A. & Gilson, M. K. Prediction of pH-dependent properties of proteins. *J. Mol. Biol.* **238**, 415–436 (1994).
51. Chan, P. & Warwicker, J. Evidence for the adaptation of protein pH-dependence to subcellular pH. *BMC Biol.* **7**, 69 (2009).
52. Dalkas, G. A., Teheux, F., Kwasigroch, J. M. & Rooman, M. Cation- $\pi$ , amino- $\pi$ ,  $\pi$ - $\pi$ , and H-bond interactions stabilize antigen-antibody interfaces. *Proteins: Struct. Funct., Bioinf.* **82**, 1734–1746 (2014).
53. Tartaglia, G. G., Pechmann, S., Dobson, C. M. & Vendruscolo, M. Life on the edge: a link between gene expression levels and aggregation rates of human proteins. *Trends Biochem. Sci.* **32**, 204–206 (2007).
54. Nichols, P. *et al.* Rational design of viscosity reducing mutants of a monoclonal antibody: hydrophobic versus electrostatic intermolecular interactions. *mAbs* **7**, 212–230 (2015).
55. Berman, H. M. *et al.* The Protein Data Bank. *Nucleic Acids Res.* **28**, 235–242 (2000).
56. Rose, A. S. *et al.* NGL Viewer: Web-based molecular graphics for large complexes. *Bioinformatics* **1**, 4 (2018).
57. Moutevelis, E. & Warwicker, J. Prediction of pKa and redox properties in the thioredoxin superfamily. *Protein Sci.* **13**, 2744–2752 (2004).
58. Warwicker, J. Improved pKa calculations through flexibility based sampling of a water-dominated interaction scheme. *Protein Sci.* **13**, 2793–2805 (2004).
59. Bate, P. & Warwicker, J. Enzyme/non-enzyme discrimination and prediction of enzyme active site location using charge-based methods. *J. Mol. Biol.* **340**, 263–276 (2004).
60. Schrödinger, L. L. C. The PyMOL molecular graphics system. Schrödinger, LLC (2010).
61. Birch, J. R. & Racher, A. J. Antibody production. *Adv. Drug Deliv. Rev.* **58**, 671–685 (2006).
62. Warwicker, J. Simplified methods for pKa and acid pH-dependent stability estimation in proteins: removing dielectric and counterion boundaries. *Protein Sci.* **8**, 418–425 (1999).
63. Beroza, P., Fredkin, D. R., Okamura, M. Y. & Feher, G. Protonation of interacting residues in a protein by a Monte Carlo method: application to lysozyme and the photosynthetic reaction center of *Rhodobacter sphaeroides*. *Proc. Natl. Acad. Sci.* **88**, 5804–5808 (1991).

## Acknowledgements

Members of the Curtis and Warwicker groups are thanked for discussion and providing feedback. The authors would like to acknowledge the assistance given by IT Services at The University of Manchester. UK EPSRC grant EP/N024796/1.

## Author Contributions

M.H. and J.W. conceived and conducted the experiment(s), analysed the results and reviewed the manuscript.

## Additional Information

**Supplementary information** accompanies this paper at <https://doi.org/10.1038/s41598-018-36950-8>.

**Competing Interests:** The authors declare no competing interests.

**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 4.0 International License, which permits use, sharing, adaptation, 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 license, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons license, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons license 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 license, visit <http://creativecommons.org/licenses/by/4.0/>.

© The Author(s) 2019