

Supervised Learning Model Predicts Protein Adsorption to Nanotubes

Nicholas Ouassil,^{†1} Rebecca L. Pinals,^{†1} Jackson Travis Del Bonis-O'Donnell,¹ Jeffrey Wang,¹ Markita P. Landry^{*1,2,3,4}

¹ Department of Chemical and Biomolecular Engineering, University of California, Berkeley, California 94720, United States

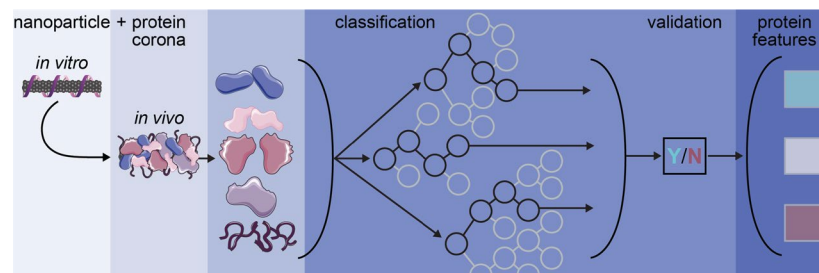
² Chan-Zuckerberg Biohub, San Francisco, California 94158, United States

³ Innovative Genomics Institute (IGI), Berkeley, California 94720, United States

⁴ California Institute for Quantitative Biosciences, QB3, University of California, Berkeley, California 94720, United States

[†] Co-authors

^{*} Corresponding author, landry@berkeley.edu



Abstract

Engineered nanoparticles are advantageous for numerous biotechnology applications, including biomolecular sensing and delivery. However, testing the compatibility and function of nanotechnologies in biological systems requires a heuristic approach, where unpredictable biofouling often prevents effective implementation. Such biofouling is the result of spontaneous protein adsorption to the nanoparticle surface, forming the “protein corona” and altering the physicochemical properties, and thus intended function, of the nanotechnology. To better apply engineered nanoparticles in biological systems, herein, we develop a random forest classifier (RFC) trained with proteomic mass spectrometry data that identifies which proteins adsorb to nanoparticles. We model proteins that populate the corona of a single-walled carbon nanotube (SWCNT)-based optical nanosensor. We optimize the classifier and characterize the classifier performance against other models. To evaluate the predictive power of our model, we then apply the classifier to rapidly identify and experimentally validate proteins with high binding affinity to SWCNTs. Using protein properties based solely on amino acid sequence, we further determine protein features associated with increased likelihood of SWCNT binding: proteins with high content of solvent-exposed glycine residues and non-secondary structure-associated amino acids. Furthermore, proteins with high leucine residue content and beta-sheet-associated amino acids are less likely to form the SWCNT protein corona. The classifier presented herein provides an important tool to undertake the otherwise intractable problem of predicting protein-nanoparticle interactions, which is needed for more rapid and effective translation of nanobiotechnologies from *in vitro* synthesis to *in vivo* use.

Introduction

Engineered nanoparticles are poised to transform how we undertake biological sensing,(1–3) imaging,(4–6) and delivery:(7–9) nanoscale materials enable localization within otherwise inaccessible biological environments and exhibit highly tunable physicochemical properties to tailor function. Different nanoparticle platforms offer application-dependent advantages, such as near-infrared fluorescent nanoparticles for through-tissue imaging(10, 11) or biodegradable nanoparticles for *in vivo* delivery.(12–14) In particular, single-walled carbon nanotubes (SWCNTs) are well-suited for biological sensing and imaging due to their tissue-transparent and photostable near-infrared fluorescence, in addition to their readily modifiable surface.(15–17) As such, SWCNTs have been functionalized with biomolecules including single-stranded DNA to create neurotransmitter nanosensors,(18–20) with peptide mimetics to form protein nanosensors,(21) and with proteins to construct viral nanosensors.(22) Similarly, the large SWCNT surface area enables cargo attachment such that SWCNTs can be loaded with DNA plasmids or small interfering RNAs, translocating these functional biomolecules into cells for gene expression and silencing applications.(23–26) Optimizing these biomolecule-nanoparticle interactions is key in enhancing nanotechnology function, and a deeper understanding of these interfacial interactions would enable more rational conjugate designs. As such, the capability to predict nano-bio interactions would aid the design phase of nanobiotechnologies by lessening the need to experimentally test innate interactions of each biomolecule with each nanoparticle of interest.

Although such aforementioned nano-bio interactions are required for function, biofouling of nanobiotechnologies results from undesired nano-bio interactions that often inhibit intended nanoparticle function. Functionalized SWCNTs and other nanotechnologies more broadly suffer from as-of-yet unpredictable interactions with the biological environments in which they are applied. When engineered nanoparticles are introduced into biological systems, endogenous proteins rapidly bind to the nanoparticle surface.(27–29) This phenomenon is known as protein corona formation. Protein adsorption often decreases the ability of the nanoparticle to interact with its surrounding environment, such as sensing nearby analytes or navigating biological barriers.(30, 31) For sensing applications, protein corona formation sterically hinders access of target analytes to the nanosensor surface and unpredictably changes the sensor baseline required for accurately calibrating and quantifying analyte levels.(32–34) For imaging and delivery applications, the protein corona modifies the *in vivo* trafficking, biodistribution, biocompatibility, and overall functionality of nanotechnologies.(35, 36) Consequently, the corona often reduces the efficiency with which nanoparticle-based contrast agents or cargo-filled vehicles reach their intended locations.(37–39) Passivation with anti-biofouling ligands such as polyethylene glycol (PEG) is a promising technique to reduce protein-binding on foreign surfaces and to retain the pristine, as-designed nanoparticle properties.(39–44) Still, the protein corona remains as a complex and poorly understood phenomenon limiting the efficacious application of nanotechnologies in biological systems. Knowledge of the proteins adsorbed in this corona phase would enable better prediction of the biological identity, and thus fate, of the applied nanotechnologies.(45, 46) Limitations in our understanding of corona formation arise from a convolution of diverse nanoparticle properties (dominated by surface characteristics) and the complexity of biological environments.(28, 41) Experimental testing to fully characterize the protein corona on all synthesized nanoparticle constructs within all intended biological environments is laborious and costly. While recent work has made headway toward high-throughput experimental methods,(47, 48) the most common strategies still rely on labor-intensive mass spectrometry-based proteomics.(41, 49) The ability to predict the protein corona that will form on nanoparticles *in vivo* remains a challenge that, if overcome, would move the field toward better clinical translatability.

Pattern recognition techniques, including those of machine learning, offer a route to characterize protein-nanoparticle interactions in a high-throughput manner across this extensive design space of nanoparticles

applied in different biological systems. Previous work pioneering this idea applied random forest classification to predict what proteins adsorb to silver nanoparticles in biologically relevant environments,(49) and has been since expanded to larger nanoparticle libraries.(50) However, certain aspects stand to be refined, such as setting the threshold of whether a protein is classified as in or out of the corona. Other work has examined protein-nanoparticle complexes using a fluorometric assay to guide prediction of corona formation, though issues arise in characterizing fluorophore interactions on graphene-based substrates.(51) More broadly, most predictive modeling efforts involving nanoparticles applied in biology consider cellular- or organism-level responses, such as cellular association,(52, 53) toxicity,(54) *in vivo* fate,(46) and therapeutic efficacy.(53, 55) Toward protein-SWCNT conjugate design, some predictive modeling has informed protein candidates that exhibit natural affinity for the graphitic SWCNT surface.(23) For example, Di Giosia *et al.* implemented a molecular docking model to determine a panel of proteins that interact with the carbon nanotube surface.(56) Yet, this strategy of predicting protein corona identity requires protein structural information and is low throughput, both computationally and in experimental validation.

Herein, we develop a classifier to predict protein-nanotube association based on physicochemical properties of proteins. Our purpose is two-fold: as one objective, we aim to predict which protein-SWCNT interactions to expect in true biological environments. This knowledge will inform implementation of anti-biofouling strategies toward effective biological application of nanoparticles. Our second objective is to predict high-affinity protein binders to SWCNTs and protein features associated with such binding affinity, to improve the process of protein-nanoparticle construct design.(23) Toward these ends, we build and validate a random forest classifier to predict protein adsorption to SWCNTs. We relate protein properties (derived from protein sequence data) to whether proteins are in or out of the corona phase on SWCNTs (experimentally determined by quantitative mass spectrometry-based proteomics). Specifically, we focus on protein corona formation on (GT)₁₅-SWCNTs due to their demonstrated applicability for dopamine sensing, however, the workflow is generalizable to other nanoparticles.(18, 19) We train our classifier using mass spectrometry-based proteomic data characterizing the corona formed on (GT)₁₅-SWCNTs in two relevant bioenvironments: the intravenous environment (blood plasma) and the brain environment (cerebrospinal fluid).(57) We find that our classifier can precisely target a small number of proteins that adsorb to our nanoparticle. Furthermore, we identify population distribution changes among the most important protein properties to gain insight on how our classifier identifies positive targets. Namely, high content of glycine residues (particularly solvent-exposed residues) and amino acids not associated with secondary structure domains (alpha helix, beta sheet, and turns) lead to favorable SWCNT binding, whereas high content of leucine residues and amino acids associated with planar beta-sheet domains lead to unfavorable SWCNT binding. Finally, we test our model with an entirely new set of proteins and perform quantitative protein adsorption experiments to validate the model's in vs. out of corona predictions.(33) Our results suggest that this classifier can serve as a valuable method to both overcome the high failure rate in translating nanotechnologies from *in vitro* validation to *in vivo* deployment, and to aid in rational design of future nano-bio tools.

Results

Experimentally determined protein corona composition on (GT)₁₅-SWCNTs

The training data was experimentally generated from a selective adsorption assay that quantifies proteins present on (GT)₁₅-SWCNT nanoparticles incubated in either human blood plasma or cerebrospinal fluid (CSF) of the brain, characterized using liquid chromatography tandem mass spectrometry (LC-MS/MS).(57) This experimental dataset reveals the corona components with quantitative protein amounts. The absolute protein abundance and relative enrichment or depletion (compared to the control sample of the biofluid alone) was used to indicate whether a particular protein was considered to be in the corona, as will be

described in a later section. We included four training datasets: (GT)₁₅-SWCNTs in either blood plasma, cerebrospinal fluid, total biofluid datasets, and total naïve. Total and total naïve only differ by one variable, namely, the former contains the biofluid phase from which the protein originated. Although we focus on protein corona characterization with one nanoparticle type, SWCNTs, it is worthwhile to note that these protein datasets do not require any information regarding the nanoparticle itself. The only location where nanoparticle data is included is the named class (i.e., in or out of the corona).

Protein property database development from protein sequence

We next curated a protein property database to use with our classifier. We used the amino acid sequence of each protein from the annotated protein database, UniProt,(58) to construct an array of predicted physicochemical protein properties with the BioPython package (**Table S1**).⁽⁵⁹⁾ Our protein property database requires access to only the amino acid sequence, enabling expansion to new proteins as needed for future experimental datasets or nanoparticle-binding proteins of interest. Although UniProt also provides biological protein properties (such as gene ontology, sequence annotations, and specific functional regions), our final classifier was based solely on amino acid sequence data to avoid potential issues of less well-studied proteins that have no empirically derived properties and/or no annotated features (classifier performance comparison in **Figure S1**). Importantly, developing a database in this manner expands the number of possible proteins that can be tested because the classifier does not require prior information on the annotated protein sequence nor interactions between the protein and nanoparticle.

The amino acid sequence of a protein provides valuable information including the percentage of a specific amino acid within the full protein; however, spatial organization is disregarded. To complement the sequence-derived dataset, we added the parameter of solvent accessibility that estimates the exposed protein surface area. We implemented NetSurfP 2.0⁽⁶⁰⁾ to predict the number of exposed residues of a particular protein using the amino acid sequence, normalized by either the total number of amino acids or the total number of exposed amino acids. These two choices of normalization provide information on amino acid content on the surface relative to the full protein or relative to only other surface-exposed residues, respectively. To collate this data, we programmatically created submissions from UniProt protein sequence entries and rapidly collected data, aligning with our goal of creating an easily expandable database.

Thresholding to determine protein placement: in or out of the corona

The decision of whether a protein was categorized as in or out of the corona was made using the protein abundance data from LC-MS/MS experiments. Proteins were placed into the corona based on two criteria: (i) relative change and (ii) an abundance threshold. First, if a protein was more abundant in the nanoparticle-bound case than it was in the control solution of the native biofluid without any nanoparticles present (i.e., enrichment on nanoparticle), then it was classified as in the corona. Second, the remaining proteins were ordered by abundance and fit to an exponential distribution. Increasing the power of the exponential leads to a higher in-corona threshold, placing fewer proteins in the corona. Importantly, this thresholding approach reflects that lower abundance of a protein in the corona relative to its abundance in the biofluid (i.e., depletion on nanoparticle) does not necessarily qualify a protein as out of the corona: a protein that is significantly depleted can still be present in the corona with a high absolute quantity. The thresholding method that we have developed is discussed further in the Methods section, with comparison to Otsu's method as a common form of thresholding applied in image analysis.

Random forest classifier development using established protein property database

In line with previously published work, we implemented a random forest classifier (RFC) to classify proteins as in or out of the corona phase on (GT)₁₅-SWCNT nanoparticles. We chose to pursue ensemble methods due to the concern of overfitting the classifier. To confirm the choice of RFC over other potential classifiers, we tested an assortment of classifier types (**Figure S2**). The highest performing classifiers were the RFC

and a Bagging classifier using decision trees based on a sum of accuracy, precision, and recall. We selected the RFC for this study because the precision (0.671) and accuracy (0.747) values were superior to that of the Bagging tree, while retaining similar area under the receiver operating curve (AUC; 0.700). AUC is a frequently used measure for understanding sensitivity and specificity of the classifier. Moreover, the RFC provided the highest precision (positive predictive value), which is favorable for the most straightforward application of classifier output for nanobiotechnology optimization. However, the Bagging tree did perform better than the RFC in recall (Bagging Tree, 0.528; RFC 0.505).

Due to the imbalance in our LC-MS/MS experimental dataset (i.e., unequal number of proteins in either class), we up-sampled our minority class (in corona; ~30% in corona without up-sampling in combined dataset). This up-sampling ensures that each time the classifier was trained, we were able to recover an appropriate amount of the minority class. For this reason, the classifier was validated using a stratified shuffle split repeated 100 times. Moreover, we noticed that generalization of this classifier could also be quite poor, especially when considering recall which was below 0.5. To address this issue, a synthetic minority over-sampling technique (SMOTE)(61) was implemented to generate new “proteins” in the minority class (in corona). This analysis revealed that the most precise and accurate results for our classifier were obtained when the minority/majority ratio in SMOTE was 0.5/1.0 (**Figure S3**; precision 0.678 and accuracy 0.749). The recall of our classifier was improved marginally to from 0.497 to 0.512. Introducing the described methods widely expanded the number of proteins that were placed in the corona, thus enhancing the classifier’s generalization ability. However, this SMOTE ratio offers a tunable handle: if an experimenter preferred higher recall values, a ratio of 1 provides a recall of 0.587, although reducing precision to 0.571.

RFC verification

Using an RFC, classification tests were run on the total naïve dataset of proteins marked as being in or out of the corona using the aforementioned thresholding method. The classifier performance was scored for a range of thresholding powers (**Figure 1a**). The classifier was then refreshed and the standard protocol for training the classifier was repeated to gather metrics related to classification: accuracy, area under the receiver operating curve (AUC), precision, and recall. The metrics were recorded until a thresholding power of 3.5, at which point higher powers considerably reduce the number of proteins counted in the plasma corona and many metrics drastically decline in their performance. We ultimately selected a power of 2.25 because this power provided the best compromise between accuracy (0.747), AUC (0.691), and precision (0.648), while only suffering slightly in recall (0.570). All reported results for the remainder of this work use a power of 2.25 for placing proteins in the nanoparticle corona.

During the development of our classifier, stratified shuffle split validation was used to check the success of our classifier regarding accuracy, area under the receiver operating curve, recall, and precision. The dataset was divided into a training and test set at the beginning of each split, then the training data was fit to an untrained classifier. Next, the test set was used to make predictions and compared with our true answers. The result of this classifier was saved and the process was repeated with the classifier naïve at the beginning of each iteration, as graphically depicted in **Figure 1b**. This method was used to ensure that the subset of proteins generated more accurate metrics for the classifier, considering each protein revolves into the testing set during one of the folds. Statistics represented in this work are generated from the n trials used in this verification step.

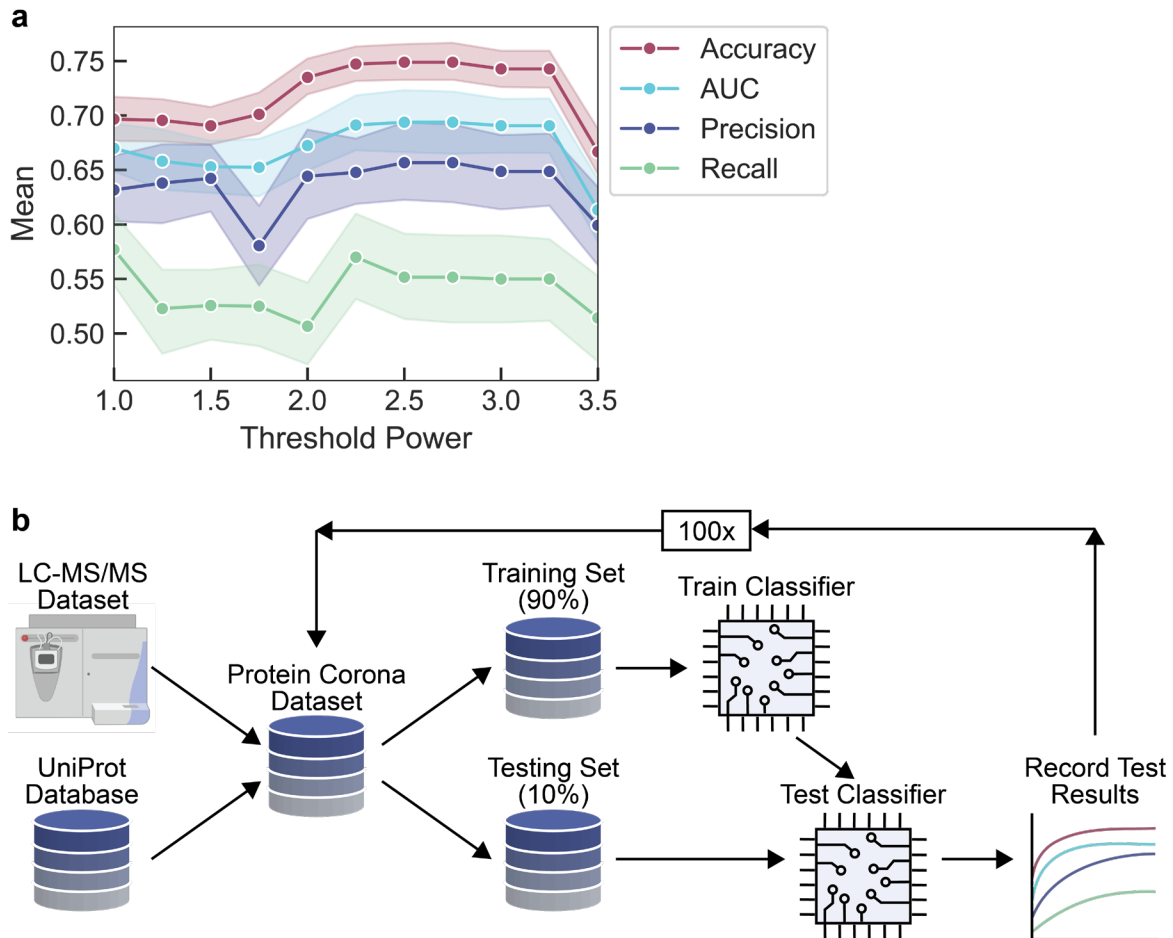


Figure 1. Random forest classifier (RFC) development and workflow for determining proteins as in vs. out of the corona phase on $(GT)_{15}$ -SWCNTs. **(a)** Metrics of accuracy, area under the receiver operating curve (AUC), precision, and recall recorded as a function of threshold power for classifying proteins as in vs. out of the corona. A threshold power value of 2.25 was selected for subsequent analyses due to the optimal combination of the recorded metrics. Shaded error bars represent 95% confidence intervals. **(b)** RFC workflow used in splitting-based predictions. Liquid chromatography-tandem mass spectrometry (LC-MS/MS) experimentally provides protein corona composition. LC-MS/MS data is combined with protein properties derived from the protein sequence (UniProt database) to form a total dataset. The total dataset is split 90% into training data and 10% into test data. Training data trains a reset classifier then test data is used to score the trained classifier. Results are recorded and the process is repeated.

Throughout this process, results were collated from each round of the classifier. The first trial was the difference of two datasets, total vs. total naïve (**Figure 2a**). The only difference between these two datasets was the inclusion of one Boolean column that dictates from which biofluid a protein originated. We observe that the inclusion of this “biofluid of origin” information does not improve the classification ability on our complete dataset. Thus, we deemed this column unnecessary to include for future runs. Moreover, keeping this column would have made our classifier less generic when selecting new proteins that may not be present in blood plasma or CSF.

We next trained the classifier on corona proteins present in one biofluid and attempted to predict corona proteins present in another biofluid. For this case, instead of splitting the training data 90%/10%, the classifier was trained on one complete dataset, then the second complete dataset was subset into a testing

set. We repeated this approach 100 times to generate statistics for the classifier. We notice similar results in AUC (CSF: 0.702, plasma: 0.691) and accuracy (CSF: 0.687, plasma: 0.706) independent of which biofluid the classifier was trained on (**Figure 2a**). However, there is a difference in precision (CSF: 0.469, plasma: 0.649) and recall (CSF: 0.676, plasma: 0.577) for each of these classifiers, arising from the inclusion of a few of proteins that are present in the corona formed on (GT)₁₅-SWCNTs from one biofluid and are not present in the corona formed on (GT)₁₅-SWCNTs from the other biofluid (e.g., serotransferrin found in the CSF corona and haptoglobin found in the plasma corona). This discrepancy occurs because our classifier has no context of which proteins are in the corona formed from which biofluid, and thus no manner of adjusting to proteins portraying contradicting adsorptive behavior across biofluids. However, this classification discrepancy only occurs for a few proteins (13 proteins out of 38 duplicate proteins, within 174 total proteins).

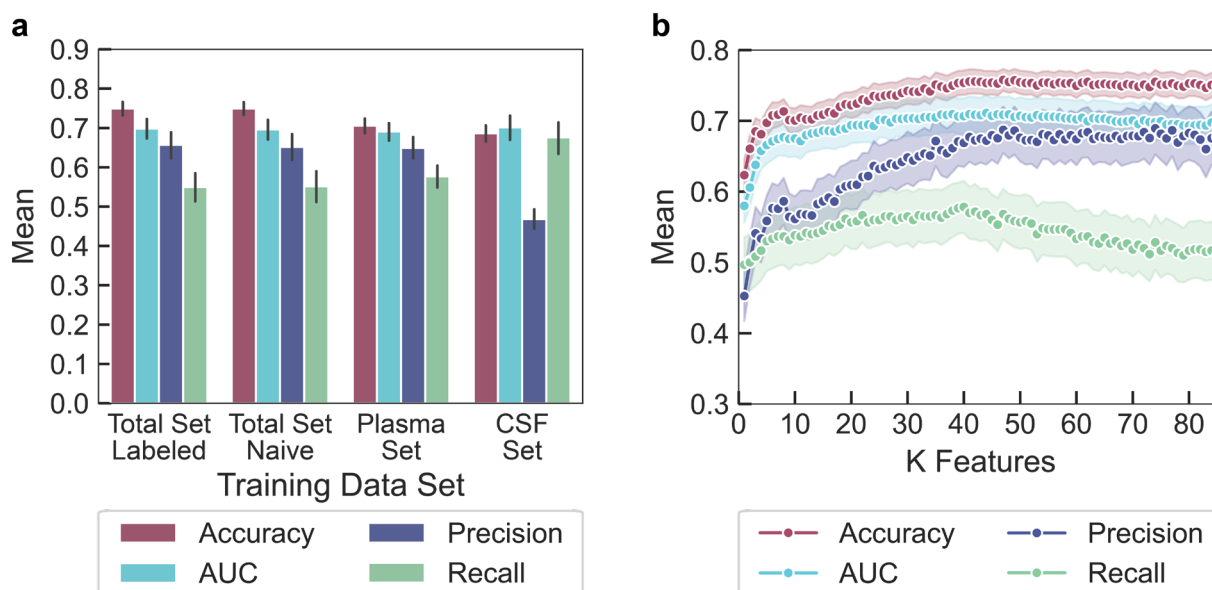


Figure 2. Classifier results on varied training datasets and with varied feature inputs. **(a)** RFC is trained on the full protein set (with or without the label of origin biofluid) or each individual biofluid (plasma or CSF). Negligible differences arise between the RFC's ability to classify the total set with or without the biofluid label (total set labeled compared to total set naïve), denoting that this biofluid label feature does not resolve the cross-fluid classification problem. Training the RFC on one biofluid and testing against the second biofluid produced similar metrics except for precision, attributable to a few proteins labeled in the corona of one biofluid but not the other. Error bars represent 95% confidence intervals. **(b)** RFC is trained on the total naïve protein corona set, with features sorted by ANOVA and added to the classifier from highest to lowest importance. At approximately 40 features, classification ability begins to plateau for all metrics except recall. By 89 features, there is a decline in recall but marginally enhanced precision. Shaded error bars represent 95% confidence intervals.

Feature analysis for importance and correlation with class predictions

During the development of our model, 89 protein features were mined as potentially important to classify these proteins as in vs. out of the nanoparticle corona (**Table S1**). Each feature was examined for the extent of contribution to the overall classification ability of the system using an ANOVA test. Features were added in one-by-one until the classifier had scored all 89 features (**Figure 2b**). This analysis indicates that there is a minimum number of features of approximately 10 to result in sufficient classification ability. We also note that use of approximately 40 features provides a maximum for recall and AUC metrics. If we include up to 89 features, we see a marginal increase in the precision ability of our classifier with a marginal

decrease in AUC and recall. As such, experimenters can tune the number of features depending on whether precision or recall is more important: more precise results will be better for experimenters using this tool to correctly identify new nanoparticle-binding proteins, while higher recall results will be better when the opportunity cost of missing a positive corona contributor is more problematic than including a false positive.

Using the feature ranking by ANOVA, the top ten protein features influencing protein adsorption to (GT)₁₅-SWCNTs were identified (**Table 1**). Since RFCs do not provide correlational information (i.e., whether a high importance ranking positively or negatively influences protein adsorption), we calculated basic kernel density estimates on distributions of these features and we examined how these distributions changed to hypothesize correlations (**Figure 3**; top ten feature distributions in **Figure S4**). We find that the fraction of solvent-exposed amino acid glycine (normalized to either the total exposed amino acid count or the total amino acid count), the fraction of amino acid glycine, and the fraction of predicted non-secondary structure-associated amino acids correlate positively with the protein being in the corona. Conversely, the fraction of amino acid leucine and the fraction of beta-sheet secondary structure-associated amino acids correlate negatively with being in the corona. Previously, we linearly regressed the log-fold change (ratio of protein amount in the corona vs. in the native biofluid) against physicochemical protein properties to understand protein features that govern corona formation, using these same nanoparticle and biofluid data sets.(57) In this experimental dataset, high leucine content was similarly determined to be less favorable for protein adsorption. High glycine content was found to be associated with more favorable protein adsorption when included in the regression analysis. However, glycine contribution was not evaluated in the original regression due to correlation with other protein features, as the calculated variance inflation factor was greater than the set threshold value.(57) As such, glycine content impact could not be deconvoluted from other protein properties. This analysis highlights a benefit of the current RFC over the previously applied linear regression approach, where co-dependent variables must be proactively excluded in the latter case. It should further be noted that secondary structure features were not included in the protein property database for the linear regression analysis due to data sparsity, a problem that is overcome with the current study by implementing BioPython to predict such features from the amino acid sequence without relying on protein structure annotations.

Our analysis of the top protein features promoting corona binding indicates that more flexible proteins are favorable to bind to (GT)₁₅-SWCNTs, as inferred by high glycine content and less strict secondary structural domains. This result is in agreement with previous experimental work demonstrating that peptides and small molecule ligands possessing more conformational flexibility bind more readily to carbon nanotubes.(62, 63) Increased adsorption propensity suggests that more flexible proteins are able to maximize favorable surface contacts with the highly curved SWCNT, in comparison to rigid proteins with energetic penalties associated with adopting new surface-adsorbed conformations. Interestingly, flexibility itself appears in the bottom ten most important protein features for protein corona formation (**Table S2**). This measure of flexibility was calculated by Vihinen *et al.* using normalized B-factors (i.e., Debye–Waller factors) for each residue. B-factors incorporate the dependence on neighboring amino acids with a 9-residue sliding window averaging approach.(64) With this method, glycine is only the top 8th most flexible residue, posited to be because glycine frequently appears on the protein surface and interior, as well as in tight turns. The restricted mobility of glycine in the interior and turn motifs may reduce the overall flexibility value. As such, our result that high glycine content specifically located on the protein surface is an enriched feature in the corona phase indicates that protein flexibility leads to higher protein corona binding on SWCNTs. In comparison to previous literature, glycine has been found to display a relatively low magnitude, yet still favorable, free energy change upon binding to pristine SWCNTs, as determined by enhanced sampling molecular dynamics.(65) However, this study was done at the scale of single amino acid analogs. Accordingly, this study disregards the full-protein structural context of each amino acid. Finally, intrinsically disordered proteins have been demonstrated to disperse SWCNTs stably in the aqueous phase even under

mild sonication conditions.(66) Although the non-structure-associated amino acid content that we report is not equivalent to intrinsically disordered domains, our result is in line with these previous experimental findings.

In contrast, our analysis of top protein features that deter corona binding reveals that proteins high in the aliphatic, hydrophobic amino acid leucine and proteins with more planar beta-sheet character are not expected to bind to (GT)₁₅-SWCNTs. The finding that hydrophobic leucine does not increase SWCNT binding is not necessarily intuitive, considering that the SWCNT surface is highly hydrophobic. However, this result recapitulates prior literature that nonspecific hydrophobic interactions alone do not drive corona binding;(62, 65, 67, 68) rather, aromatic, hydrophobic amino acids, especially tryptophan, are repeatedly found to be the highest binders to SWCNTs.(62, 67–70) For physical context, (GT)₁₅ ssDNA is observed to wrap helically around SWCNTs based on both experiment(71, 72) and modeling,(73, 74) though only covering ~2-25% of the aromatic SWCNT surface.(73, 75–77) Interestingly, the RFC did not highlight aromatic amino acid content (tryptophan, tyrosine, or phenylalanine) as top features for corona binding, although the fraction of exposed tryptophan is the fifth most favorable feature. In studies of isolated amino acids or short peptide sequences, aromatic amino acids seemingly drive adsorption to SWCNTs via π - π interactions with the SWCNT surface. However, in the full protein context, these π - π interactions may not be sufficient to drive initial protein contact with the SWCNT surface, as these hydrophobic amino acids are expected to be predominately buried in the folded protein core. Finally, the finding that high content of amino acids associated with beta-sheet structures leads to low protein adsorption to SWCNTs indicates the difficulty for planar protein secondary structures to adapt to the highly curved nanoparticle surface. This result is in line with previous work demonstrating that the extremely high curvature of carbon nanotubes must be aligned at the amino acid level of proteins, much less the secondary structure level:(62, 67) typical amino acid side chain lengths are on the order of 0.1-0.5 nm, in comparison to the SWCNT diameter of approximately 1 nm. Overall, the identification of these features is important in helping to predict high biofouling protein types or rationally selecting proteins to bind to nanoparticles prior to testing them experimentally.

Table 1. Ordered importance of protein features by ANOVA.

Ranking	Feature
1	% Amino acid - leucine
2	% Exposed relative to total exposed amino acids - glycine
3	% Secondary structure-associated amino acids - non-structure associated
4	% Exposed relative to total amino acids - glycine
5	% Amino acid - glycine
6	% Secondary structure-associated amino acids - sheet
7	% Exposed relative to total amino acids – tryptophan
8	% Exposed relative to total amino acids – histidine
9	% Exposed relative to total exposed amino acids - alanine
10	% Exposed relative to total exposed amino acids - tryptophan

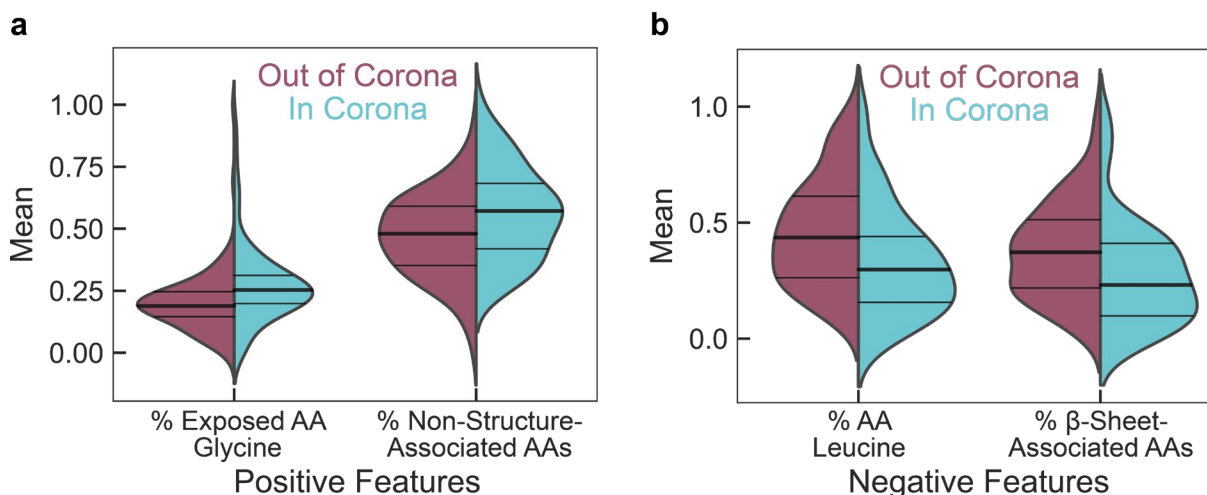


Figure 3. Distribution of the top four normalized feature values for proteins characterized as out of the corona phase (red) vs. in the corona phase (blue) on (GT)₁₅-SWCNTs. Protein features that (a) positively influence or (b) negatively influence the probability of a protein being classified as in the corona are denoted by distribution shifts toward 1 or 0, respectively. (a) Positive features include (left) the fraction of solvent-exposed amino acid (AA), glycine, relative to only the solvent-exposed amino acids and (right) the fraction of amino acids not associated with any specific secondary structure motifs. (b) Negative features include (left) the fraction of amino acid, leucine, and (right) the fraction of amino acids associated with a beta-sheet secondary structure.

Experimental validation of protein binding to SWCNTs

To test the predictive value of our supervised learning model, we applied our classifier to rank new nanotube-binding proteins and next experimentally tested the expected protein binding order. The classifier was used to predict interaction affinity of over 2,000 total proteins (available for batch download through the UniProt database(58)) with (GT)₁₅-SWCNT nanoparticles. Importantly, these proteins represent a broad class of functions and sub-cellular locations, and are distinct from those present in the plasma and CSF training datasets. Protein binding propensity was determined with associated binding probabilities, as summarized in **Table S3**. We then implemented a corona exchange assay to measure real-time, in-solution protein binding dynamics on the nanotube surface, as described previously.(33) Briefly, the ssDNA originally adsorbed on the SWCNT surface is fluorescently labeled with a Cy5 fluorophore. When near the SWCNT surface, the fluorophore is in a quenched state. Upon addition of proteins, proteins will differentially bind to the SWCNT and cause various degrees of ssDNA desorption, as denoted by de-quenching of the Cy5 fluorophore. Thus, fluorescence tracking of the Cy5-ssDNA provides a proxy for protein binding on the SWCNT without requiring fluorescent labeling or other modification of the protein.

The corona exchange assay was used to test a panel of proteins predicted to be in the corona (probability > 0.5) vs. out of the corona (probability < 0.5). Specifically, we tested the protein panel: CD44 antigen and TAR DNA-binding protein 43 (TDP-43) (predicted to adsorb to (GT)₁₅-SWCNTs) and transgelin, lysozyme C, ribonuclease pancreatic (RNase A), syntenin-1, L-lactate dehydrogenase A chain (LDH-A), and glutathione S-transferase (GST) (predicted to not adsorb to (GT)₁₅-SWCNTs) (classifier results listed in **Table S3**). Protein adsorption based on the end-state fluorescence values matched classifier predicted outcomes of in vs. out of the corona: addition of CD44 antigen and TDP-43 both resulted in sizeable ssDNA desorption from the SWCNT surface, whereas all proteins predicted to be out of the corona produced less ssDNA desorption (**Figure 4a**). However, deviations from exact orderings of predicted outcomes arise within both groups of proteins. For example, the relative ordering of CD44 antigen as the top binding protein followed by TDP-43 is reversed. However, the predicted in-corona probabilities of these two proteins differs

by less than 2%. To provide a metric of predicted vs. measured monotonicity, the Spearman's rank-order correlation coefficient was calculated to be 0.619 in comparison with 0.750 for a previous protein panel comparing DNA desorption end-state vs. proteomic mass spectrometry-derived end-state (**Figure 4c**).⁽⁵⁷⁾ Predicted protein binding probabilities were also compared to rate constants fit to the ssDNA desorption dynamics from the SWCNT surface for each injected protein (kinetic model and fits in SI, **Figure S5**). It is expected that more protein binding would correlate with a larger ssDNA desorption rate constant. However, there is poor correlation between the RFC-predicted end-state and experimental dynamics of protein-SWCNT interactions, which may be reconciled with the fact that the RFC was trained on the end-state protein corona rather than the corona composition at earlier time points.

Experimental validation was repeated for the protein panel with Cy5-(GT)₆-SWCNTs, as the shorter ssDNA oligomer is displaced more readily and thus displays a greater spread in desorption rates and values between protein species (**Figure 4b**). The resultant protein panel binding order was largely the same as that of Cy5-(GT)₁₅-SWCNTs, with a slightly higher Spearman's correlation coefficient of 0.667 (**Figure 4d**). These results confirm that the protein binding observed experimentally is mainly driven by the protein interacting directly with the SWCNT nanoparticle surface. Comparison of fit rate constants vs. predicted in-corona probabilities reveals a better correlation than that of (GT)₁₅, with the exception of RNase A (**Figure S5d**).

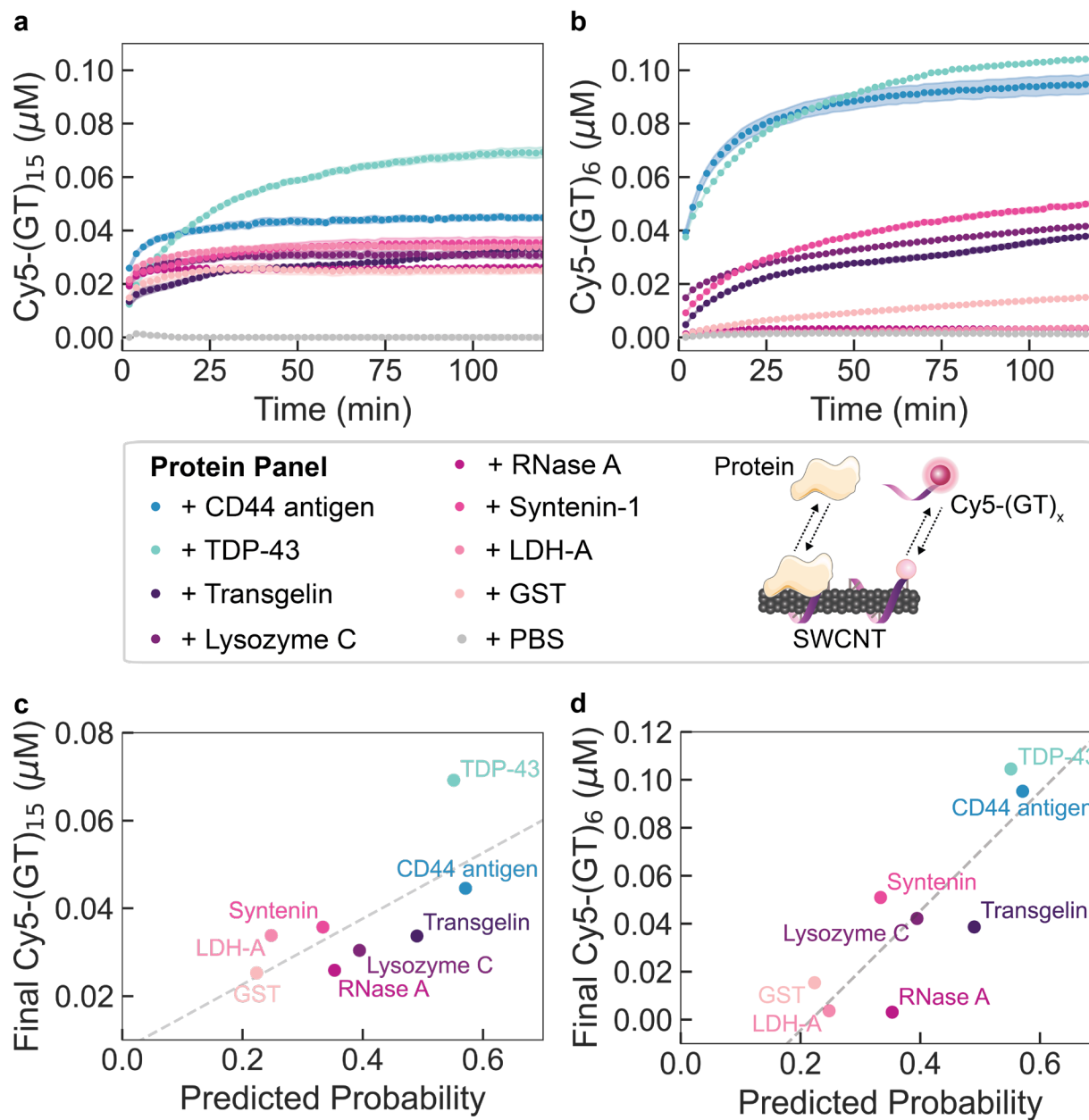


Figure 4. Protein corona dynamics assessed for binding of predicted proteins to (GT)_x-SWCNTs. **(a-b)** A corona exchange assay determines binding of a protein panel (each at 80 mg L⁻¹ final concentration) to **(a)** (GT)₁₅-SWCNTs or **(b)** (GT)₆-SWCNTs (each at 5 mg L⁻¹ final concentration). ssDNA desorption from the SWCNT serves as a proxy for protein adsorption. Proteins are predicted by the RFC to be in the corona (probability > 0.5; blue-green colors) or out of the corona (probability < 0.5; purple-pink colors). The protein panel includes: CD44 antigen and TAR DNA-binding protein 43 (TDP-43) (predicted to be in the corona) and transgelin, lysozyme C, ribonuclease pancreatic (RNase A), syntenin-1, L-lactate dehydrogenase A chain (LDH-A), and glutathione S-transferase (GST) (predicted to be out of the corona). Phosphate-buffered saline (PBS) is injected as a control and desorbed ssDNA is normalized to this initial value. Shaded error bars represent standard error between experimental replicates (N = 3). **(c-d)** End-state desorbed ssDNA is compared to the RFC predicted in-corona probability for **(c)** (GT)₁₅-SWCNTs and **(d)** (GT)₆-SWCNTs.

Examining the protein identities, it is interesting to note that lysozyme has previously been demonstrated to strongly interact with and disperse pristine carbon nanotubes, whereby hydrophobic aromatic amino acids (tryptophan and tyrosine) and cationic amino acids (arginine and lysine) are hypothesized to drive adsorption.(78–82) Yet, here we find that lysozyme interacts less with pre-dispersed ssDNA-SWCNTs based on the corona exchange results. Therefore, strong lysozyme-SWCNT interaction may hinge upon energetic input employed during the initial SWCNT dispersion process, which likely denatures lysozyme to expose more aromatic residues. Another protein of note is CD44, which is overexpressed in cancerous states including upregulation in cancer stem cells.(83) The innate affinity for CD44 to the SWCNT surface could be applied to construct a CD44-cell targeted nanotube delivery system. These results are important in suggesting that some proteins can only be adsorbed to SWCNT nanoparticles in a partially or fully denatured state, likely compromising their enzymatic activities or protein functions.

Conclusions

In sum, we developed a classifier to predict protein adsorption on ssDNA-functionalized SWCNTs with 75% accuracy and 68% precision. Ensemble methods performed better in the corona classification task and a random forest classifier scheme was ultimately chosen and optimized. We expand upon prior predictive protein corona work by (i) leveraging quantitative protein corona data,(57) (ii) redefining corona thresholding, with corresponding prediction probabilities, (iii) establishing a method for classifying proteins based solely off of the amino acid sequence of the protein of interest, and (iv) experimentally confirming adsorption in real-time, solution phase with unmodified proteins.(33) We find that no single nor small group of protein physicochemical features best determine placement in the corona. Rather, over 40 features are useful for protein classification when optimizing all four metrics of accuracy, AUC, precision, and recall. We confirm the need for these protein features by staging them into the classifier feature-by-feature and revalidating our model. Using kernel density estimates, we elucidate protein feature correlation with proteins binding or not binding to SWCNTs. Interestingly, proteins with high solvent-exposed glycine content and more non-structure-associated amino acid content (serving as proxies for protein flexibility) are found to bind in the SWCNT corona, while proteins with high leucine content and beta-sheet-associated amino acid content are not. The classifier then enabled rapid determination of proteins predicted to enter the corona phase from a new protein set, as validated experimentally with a corona exchange assay. The use of our machine learning algorithm allows us to quickly parse protein properties from a publicly available database to determine protein features of interest for corona formation, in turn informing heuristics to rationally select proteins for nanoparticle complexation in the future, or to predict biofouling of nanotechnologies.

Our supervised learning model uses amino acid sequence-based prediction of protein corona formation, which could be generalizable across a wide range of nanoparticles and bioenvironments. *In silico* protein corona prediction will ensure that nanotechnologies can be more seamlessly implemented in biological systems without the need for experimental mass spectrometry-based proteomic characterization and analysis. In the extension of this work, nanoparticle features may be included to further enhance classification ability on different nanoparticles. However, such nanoparticle features should be readily accessible to retain the triviality of classifying new systems. Recent advances in prediction of protein properties from protein sequences alone are promising toward refinement of the protein database we have curated for this classifier, enabling inclusion of biological and annotated sequence-based protein properties that are not reliant on experimental study.(84) The ability to predict adsorption of specific proteins will enable connection to downstream cellular responses, toxicity outcomes, and overall nanotechnology functionality. The developed classifier provides a tool for both predicting key proteins expected to take part in *in vivo* biofouling and rapid prescreening of protein candidates in rationally designed nanobiotechnologies.

Methods

Database development

Protein information was downloaded from UniProt,(58) including amino acid sequences and sequence annotations. Amino acid sequences were used to generate a series of physicochemical protein properties using BioPython's Protein Analysis module **Table S1**.(59) Amino acid sequences were additionally analyzed by NetSurfP 2.0(60) to determine solvent accessibility, including relative solvent accessibility (RSA), absolute solvent accessibility (ASA), and fractions of each amino acid exposed surface area relative to either all amino acids or only other exposed amino acid surface area. The resulting data was processed and merged with the BioPython analysis. The complete database was run normalized with a Min Max Scalar from Scikit-Learn(85) before being subset and fit to the classification model. This database development method was chosen to enable facile expansion with new protein datasets. Code for this and all subsequent sections can be found in the GitHub link provided.

Criteria for in-corona placement

Using the method described previously for protein corona studies by LC-MS/MS,(57) data was obtained for proteins adsorbing to (GT)₁₅-SWCNTs in two different human biofluids: blood plasma and cerebrospinal fluid. First, proteins with abundances (A_{corona}) greater than the control of protein abundances in biofluids alone ($A_{biofluid}$) were assigned as in the corona (i.e., enriched in the corona relative to the biofluid). Second, an exponential decay, $n = n_0 \exp(-kA)$, was fit to the distribution of abundances for the remaining proteins, where n_0 and k are fitting parameters. An abundance threshold ($A_{threshold}$) was selected at a value where the exponential decay fell to a value of $n_0 \exp(-p)$, or $A_{threshold} = p/k$, where p was an optimization parameter. Proteins with an abundance greater than $A_{threshold}$ were assigned as being in the corona. We varied p between 0 and 3.5 and chose the value 2.25, which optimized the performance of the classifier following training (**Figure 1a**) and was used for the remainder of the analysis. Corona thresholding was originally completed with Otsu's method, a technique generally implemented for image thresholding.(86) However, employing Otsu's method resulted in only 3-5 proteins placed in the corona for each biofluid. Although the classifier was highly accurate at identifying these proteins, the number of proteins selected was not fully representative of the corona and we accordingly implemented our modified thresholding method described above.

Classifier selection

The use of a random forest classifier (RFC), logistic regression, bagging classifier, gradient boosting classifier, AdaBoost classifier, and XGBoost classifier were evaluated. The RFC, logistic regression, bagging classifier, gradient boosting classifier, and AdaBoost classifier were imported from Scikit-Learn.(85) The XGBoost classifier was imported from XGBoost(87) for use with Scikit-Learn. AdaBoost and bagging classifiers were tested with an underlying support vector machine, decision tree, and logistic regression. The gradient boosting classifier was tested with an underlying decision tree. XGBoost was tested with an underlying decision tree as well as 100 parallel trees.

As expected from previous literature, better performance was demonstrated with the RFC and this classifier was accordingly chosen for the remainder of the work. The classifier was next validated using a stratified shuffle split (100 repeats) validation to ensure high levels of the minority class. The minority class here is the in-corona class which has less proteins than the out-of-corona class. The shuffle split retained 10% of the dataset for corona validation. Results were collected for each fold.

Hyperparameter tuning

Using Scikit-Learn's GridSearchCV,(85) a wide range of hyperparameters, such as number or depth of trees, were tested with the classifier. With each set of hyperparameters the model was validated using the

method dictated in the **Classifier Selection** section and scored. The classifier was chosen with the hyperparameters optimized for precision using GridSearchCV.

Dimensionality reduction

To understand the effects of each feature (i.e. variable describing the protein of interest) on the total system, features were ranked using Scikit-Learn's SelectKBest function.(85) Using the ranking established from SelectKBest, the database features were unmasked one-by-one running the classifier as described in the **Classifier Selection** section until all features had been added in. Metric results were saved, and statistics were calculated.

New prediction targets

The classifier was tested against a list of 996 cytoplasmic proteins and 999 nuclear proteins (available for batch download through the UniProt database(58)), together with 45 readily accessible proteins or proteins of interest for SWCNT-based sensing and delivery applications. Amino acid sequences for these proteins were downloaded from UniProt and processed through the database development workflow described above. This new complete protein database was then processed through the classifier $k+1$ times. The first k times were completed through the described k -fold validation using the combined datasets for (GT)₁₅-SWCNTs in plasma and cerebrospinal fluid as the training and verification data. Predictions were recorded at the end of each fold. The last time new proteins were run, all available data was used to train the classifier; this last classifier then provided predictions on the new proteins, as listed in **Table S3**.

Synthesis of ssDNA-SWCNTs

Suspensions of single-walled carbon nanotube (SWCNTs) with fluorophore-labeled single-stranded DNA (Cy5-(GT)₁₅ or Cy5-(GT)₆) were prepared with 0.2 mg of mixed-chirality SWCNTs (small diameter HiPco™ SWCNTs, NanolIntegris) and 20 μ M of ssDNA (3' Cy5-labeled custom ssDNA oligos with HPLC purification, Integrated DNA Technologies, Inc.; excitation 648 nm, emission 668 nm) added in 1 mL total volume of 0.1X phosphate-buffered saline (PBS; note 1X is 137 mM NaCl, 2.7 mM KCl, 10 mM Na₂HPO₄, 1.8 mM KH₂PO₄).(33) This mixture was probe-tip sonicated for 10 min in an ice bath (3 mm probe tip at 50% amplitude, 5-6 W, Cole-Parmer Ultrasonic Processor). Cy5-ssDNA-SWCNT suspensions were centrifuged to pellet insoluble SWCNT bundles and contaminants (16.1 krcf, 30 min). The supernatant containing product was collected and Cy5-ssDNA-SWCNT concentration was calculated with measured sample absorbance at 910 nm (NanoDrop One, Thermo Scientific) and the empirical extinction coefficient, $\epsilon_{910nm}=0.02554 \text{ L mg}^{-1} \text{ cm}^{-1}$.(88) Cy5-ssDNA-SWCNTs were stored at 4°C until use, at which point the solution was diluted to a working concentration of 10 mg L⁻¹ in 1X PBS \geq 2 h prior to use.

Preparation of proteins

Proteins were sourced as listed in **Table 2**. Lyophilized proteins were reconstituted to the listed concentration in PBS, tilting intermittently to dissolve for 15 min, and filtering with 0.2 μ m syringe filters (cellulose acetate membrane, VWR International). All proteins were purified with desalting columns (Zeba Spin Desalting Columns, 0.5 mL with 7 kDa MWCO, Thermo Fisher Scientific) by washing with PBS three times (centrifuging 1500 rcf, 1 min), centrifuging with sample (1500 rcf, 2 min), and retaining sample in flow-through solution. Resulting protein concentration was measured with the Qubit Protein Assay (Thermo Fisher Scientific).

Corona exchange assay

Corona dynamics were measured as described previously.(33) Briefly, equal volumes (25 μ L) of ssDNA-Cy5-SWCNT and FAM-protein at 2X working concentration were added via multichannel pipette into a 96-well PCR plate (Bio-Rad) and mixed by pipetting. The PCR plate was sealed with an optically transparent adhesive seal (Bio-Rad) and briefly spun down on a benchtop centrifuge. Fluorescence was measured as

a function of time using a Bio-Rad CFX96 Real Time qPCR System, scanning all manufacturer set color channels (FAM, HEX, Texas Red, Cy5, Quasar 705) every 30 s at 22.5 °C, with lid heating off. Fluorescence time series were analyzed without default background correction.

Table 2. Purchased protein specifications.

Protein	Manufacturer	Catalog #	Lot #	Source	Notes
CD44 antigen	Acro Biosystems	CD4-H5226	616-784F1-G8	Human, expressed in HEK293	6X His tag; >95% purity
TAR DNA-binding protein 43 (TDP-43)	R&D Systems	AP-190	22675420A	Recombinant human, expressed in E. coli	>85% purity
Transgelin (TAGLN)	MyBioSource	MBS144070	1011PTAGLN30	Recombinant human, expressed in E. coli	20X His tag; >85% purity
Lysozyme C	Sigma	L2879	SLCF2361	From chicken egg white	≥80% purity
Ribonuclease pancreatic (RNase A)	New England BioLabs	T3018L		Purified from cow pancreas	
Syntenin-1	Novus Biologicals	NBP1-50893	1082301	Recombinant human, expressed in E. coli	6X His tag; >90% purity
L-lactate dehydrogenase A chain (LDH-A)	Sigma-Aldrich	10127230001	42032824	From rabbit muscle	
Glutathione S-transferase (GST)	Abcam	ab86775	GR3377596-1	Recombinant mouse, expressed in E. coli	>95% purity

Acknowledgements

We acknowledge support of the IGI LGR ERA, GlaxoSmithKline, and Citris/Banatao Seed Funding. We acknowledge support of a Burroughs Wellcome Fund Career Award at the Scientific Interface (CASI) (to M.P.L.), a Dreyfus foundation award (to M.P.L.), a Stanley Fahn PDF Junior Faculty Grant with Award # PF-JFA-1760 (to M.P.L.), a Beckman Foundation Young Investigator Award (to M.P.L.), an NIH MIRA award (to M.P.L.), an NSF CAREER award (to M.P.L), an NSF CBET award (to M.P.L.), an NSF CGEM award (to M.P.L.), a FFAR Young Investigator award (to M.P.L.), a CZI investigator award (to M.P.L), a Sloan Foundation Award (to M.P.L.), a USDA BBT EAGER award (to M.P.L), a USDA NIFA Award (to M.P.L), a Moore Foundation Award (to M.P.L.), a Cisco Research Center grant (to M.P.L), and a DARPA Young Investigator Award (to M.P.L.). M.P.L. is a Chan Zuckerberg Biohub investigator, a Hellen Wills Neuroscience Institute Investigator, and an IGI Investigator. N.O., R.L.P., and J.W. acknowledge the support of NSF Graduate Research Fellowships (NSF DGE 1752814). J.T.D.B.-O. acknowledges the support of an Early Investigator Research Award from the Congressionally Directed Medical Research Program through the U.S. Department of Defense. We would like to acknowledge the use of medical clipart from Servier Medical Art by Servier (<http://smart.servier.com>), licensed under a Creative Commons Attribution 3.0 Unported License.

References

1. L. Gloag, M. Mehdipour, D. Chen, R. D. Tilley, J. J. Gooding, Advances in the Application of Magnetic Nanoparticles for Sensing. *Advanced Materials* **31**, 1904385 (2019).
2. A. B. Taylor, P. Zijlstra, Single-Molecule Plasmon Sensing: Current Status and Future Prospects. *ACS Sens.* **2**, 1103–1122 (2017).
3. P. D. Howes, R. Chandrawati, M. M. Stevens, Colloidal nanoparticles as advanced biological sensors. *Science* **346** (2014).
4. J. T. Del Bonis-O'Donnell, L. Chio, G. F. Dorlhiac, I. R. McFarlane, M. P. Landry, Advances in nanomaterials for brain microscopy. *Nano Res.* **11**, 5144–5172 (2018).
5. B. R. Smith, S. S. Gambhir, Nanomaterials for In Vivo Imaging. *Chem. Rev.* **117**, 901–986 (2017).
6. G. Hong, A. L. Antaris, H. Dai, Near-infrared fluorophores for biomedical imaging. *Nature Biomedical Engineering* **1**, 1–22 (2017).
7. M. J. Mitchell, *et al.*, Engineering precision nanoparticles for drug delivery. *Nature Reviews Drug Discovery*, 1–24 (2020).
8. J. W. Wang, *et al.*, Nanoparticle-Mediated Genetic Engineering of Plants. *Molecular Plant* **12**, 1037–1040 (2019).
9. W. Poon, B. R. Kingston, B. Ouyang, W. Ngo, W. C. W. Chan, A framework for designing delivery systems. *Nature Nanotechnology* **15**, 819–829 (2020).
10. G. Hong, *et al.*, Through-skull fluorescence imaging of the brain in a new near-infrared window. *Nature Photonics* **8**, 723–730 (2014).
11. O. T. Bruns, *et al.*, Next-generation in vivo optical imaging with short-wave infrared quantum dots. *Nature Biomedical Engineering* **1**, 1–11 (2017).
12. H. Safari, *et al.*, Biodegradable, bile salt microparticles for localized fat dissolution. *Science Advances* **6**, eabd8019 (2020).
13. R. L. Ball, K. A. Hajj, J. Vizelman, P. Bajaj, K. A. Whitehead, Lipid Nanoparticle Formulations for Enhanced Co-delivery of siRNA and mRNA. *Nano Lett.* **18**, 3814–3822 (2018).
14. L. Xiao, G. Lu, Q. Lu, D. L. Kaplan, Direct Formation of Silk Nanoparticles for Drug Delivery. *ACS Biomater. Sci. Eng.* **2**, 2050–2057 (2016).
15. A. J. Gillen, A. A. Boghossian, Non-covalent Methods of Engineering Optical Sensors Based on Single-Walled Carbon Nanotubes. *Front. Chem.* **7** (2019).
16. A. A. Boghossian, *et al.*, Near-Infrared Fluorescent Sensors based on Single-Walled Carbon Nanotubes for Life Sciences Applications. *ChemSusChem* **4**, 848–863 (2011).
17. Z. Liu, S. Tabakman, K. Welsher, H. Dai, Carbon nanotubes in biology and medicine: In vitro and in vivo detection, imaging and drug delivery. *Nano Res.* **2**, 85–120 (2009).
18. S. Kruss, *et al.*, Neurotransmitter Detection Using Corona Phase Molecular Recognition on Fluorescent Single-Walled Carbon Nanotube Sensors. *Journal of the American Chemical Society* **136**, 713–724 (2014).
19. A. G. Beyene, *et al.*, Imaging striatal dopamine release using a nongenetically encoded near infrared fluorescent catecholamine nanosensor. *Science Advances* **5**, eaaw3108 (2019).
20. S. Jeong, *et al.*, High-throughput evolution of near-infrared serotonin nanosensors. *Science Advances* **5**, eaay3771 (2019).
21. L. Chio, *et al.*, Electrostatic Assemblies of Single-Walled Carbon Nanotubes and Sequence-Tunable Peptoid Polymers Detect a Lectin Protein and Its Target Sugars. *Nano Lett.* **19**, 7563–7572 (2019).
22. R. L. Pinals, *et al.*, Rapid SARS-CoV-2 Spike Protein Detection by Carbon Nanotube-Based Near-Infrared Nanosensors. *Nano Lett.* (2021) <https://doi.org/10.1021/acs.nanolett.1c00118> (February 28, 2021).

23. A. Antonucci, J. Kupis-Rozmysłowicz, A. A. Boghossian, Noncovalent Protein and Peptide Functionalization of Single-Walled Carbon Nanotubes for Biodelivery and Optical Sensing Applications. *ACS Appl. Mater. Interfaces* **9**, 11321–11331 (2017).
24. P. D. Boyer, *et al.*, Delivering Single-Walled Carbon Nanotubes to the Nucleus Using Engineered Nuclear Protein Domains. *ACS Appl. Mater. Interfaces* **8**, 3524–3534 (2016).
25. G. S. Demirer, *et al.*, High aspect ratio nanomaterials enable delivery of functional genetic material without DNA integration in mature plants. *Nat. Nanotechnol.* **14**, 456–464 (2019).
26. G. S. Demirer, *et al.*, Carbon nanocarriers deliver siRNA to intact plant cells for efficient gene knockdown. *Science Advances* **6**, eaaz0495 (2020).
27. M. P. Monopoli, C. Åberg, A. Salvati, K. A. Dawson, Biomolecular coronas provide the biological identity of nanosized materials. *Nature Nanotechnology* **7**, 779–786 (2012).
28. A. E. Nel, *et al.*, Understanding biophysicochemical interactions at the nano–bio interface. *Nature Materials* **8**, 543–557 (2009).
29. P. C. Ke, S. Lin, W. J. Parak, T. P. Davis, F. Caruso, A Decade of the Protein Corona. *ACS Nano* **11**, 11773–11776 (2017).
30. S. Tenzer, *et al.*, Rapid formation of plasma protein corona critically affects nanoparticle pathophysiology. *Nature Nanotechnology* **8**, 772–781 (2013).
31. J. S. Gebauer, *et al.*, Impact of the Nanoparticle–Protein Corona on Colloidal Stability and Protein Structure. *Langmuir* **28**, 9673–9679 (2012).
32. C. Jiang, *et al.*, Antifouling Strategies for Selective In Vitro and In Vivo Sensing. *Chem. Rev.* **120**, 3852–3889 (2020).
33. R. L. Pinals, D. Yang, A. Lui, W. Cao, M. P. Landry, Corona Exchange Dynamics on Carbon Nanotubes by Multiplexed Fluorescence Monitoring. *J. Am. Chem. Soc.* **142**, 1254–1264 (2020).
34. M. Gravely, M. M. Safaee, D. Roxbury, Biomolecular Functionalization of a Nanomaterial To Control Stability and Retention within Live Cells. *Nano Lett.* **19**, 6203–6212 (2019).
35. K. Cai, A. Z. Wang, L. Yin, J. Cheng, Bio-nano interface: The impact of biological environment on nanomaterials and their delivery properties. *Journal of Controlled Release* (2017) <https://doi.org/10.1016/j.jconrel.2016.11.034> (September 17, 2017).
36. F. Giulimondi, *et al.*, Interplay of protein corona and immune cells controls blood residency of liposomes. *Nature Communications* **10**, 3686 (2019).
37. P. S. R. Naidu, *et al.*, Elucidating the Inability of Functionalized Nanoparticles to Cross the Blood–Brain Barrier and Target Specific Cells in Vivo. *ACS Appl. Mater. Interfaces* **11**, 22085–22095 (2019).
38. Q. Dai, *et al.*, Particle Targeting in Complex Biological Media. *Advanced Healthcare Materials* **7**, 1700575 (2018).
39. N. Bertrand, *et al.*, Mechanistic understanding of in vivo protein corona formation on polymeric nanoparticles and impact on pharmacokinetics. *Nat Commun* **8**, 1–8 (2017).
40. A. L. Klibanov, K. Maruyama, V. P. Torchilin, L. Huang, Amphipathic polyethyleneglycols effectively prolong the circulation time of liposomes. *FEBS Letters* **268**, 235–237 (1990).
41. R. L. Pinals, L. Chio, F. Ledesma, M. P. Landry, Engineering at the nano-bio interface: harnessing the protein corona towards nanoparticle design and function. *Analyst* **145**, 5090–5112 (2020).
42. D. Yang, S. J. Yang, J. T. Del Bonis-O'Donnell, R. L. Pinals, M. P. Landry, Mitigation of Carbon Nanotube Neurosensor Induced Transcriptomic and Morphological Changes in Mouse Microglia with Surface Passivation. *ACS Nano* **14**, 13794–13805 (2020).
43. E. Ostuni, R. G. Chapman, R. E. Holmlin, S. Takayama, G. M. Whitesides, A Survey of Structure–Property Relationships of Surfaces that Resist the Adsorption of Protein. *Langmuir* **17**, 5605–5620 (2001).
44. Q. Wei, *et al.*, Protein Interactions with Polymer Coatings and Biomaterials. *Angewandte Chemie International Edition* **53**, 8004–8031 (2014).

45. K. A. Dawson, Y. Yan, Current understanding of biological identity at the nanoscale and future prospects. *Nature Nanotechnology* **16**, 229–242 (2021).
46. J. Lazarovits, *et al.*, Supervised Learning and Mass Spectrometry Predicts the in Vivo Fate of Nanomaterials. *ACS Nano* **13**, 8023–8034 (2019).
47. K. M. Poulsen, T. Pho, J. A. Champion, C. K. Payne, Automation and low-cost proteomics for characterization of the protein corona: experimental methods for big data. *Anal Bioanal Chem* **412**, 6543–6551 (2020).
48. R. Oliverio, *et al.*, Versatile and High-Throughput Strategy for the Quantification of Proteins Bound to Nanoparticles. *ACS Appl. Nano Mater.* **3**, 10497–10507 (2020).
49. M. R. Findlay, D. N. Freitas, M. Mobed-Miremadi, K. E. Wheeler, Machine learning provides predictive analysis into silver nanoparticle protein corona formation from physicochemical properties. *Environmental Science: Nano* **5**, 64–71 (2018).
50. Z. Ban, *et al.*, Machine learning predicts the functional composition of the protein corona and the cellular recognition of nanoparticles. *PNAS* **117**, 10492–10499 (2020).
51. Y. Duan, *et al.*, Prediction of protein corona on nanomaterials by machine learning using novel descriptors. *NanoImpact* **17**, 100207 (2020).
52. C. D. Walkey, *et al.*, Protein Corona Fingerprinting Predicts the Cellular Interaction of Gold and Silver Nanoparticles. *ACS Nano* **8**, 2439–2455 (2014).
53. D. Fourches, *et al.*, Quantitative Nanostructure–Activity Relationship Modeling. *ACS Nano* **4**, 5703–5712 (2010).
54. X. Bai, *et al.*, Toward a systematic exploration of nano-bio interactions. *Toxicology and Applied Pharmacology* **323**, 66–73 (2017).
55. G. Yamankurt, *et al.*, Exploration of the nanomedicine-design space with high-throughput screening and machine learning. *Nature Biomedical Engineering* **3**, 318–327 (2019).
56. M. Di Giosia, *et al.*, High-throughput virtual screening to rationally design protein - Carbon nanotube interactions. Identification and preparation of stable water dispersions of protein - Carbon nanotube hybrids and efficient design of new functional materials. *Carbon* **147**, 70–82 (2019).
57. R. L. Pinals, *et al.*, Quantitative Protein Corona Composition and Dynamics on Carbon Nanotubes in Biological Environments. *Angewandte Chemie International Edition* **59**, 23668–23677 (2020).
58. The UniProt Consortium, UniProt: the universal protein knowledgebase in 2021. *Nucleic Acids Research* **49**, D480–D489 (2021).
59. P. J. A. Cock, *et al.*, Biopython: freely available Python tools for computational molecular biology and bioinformatics. *Bioinformatics* **25**, 1422–1423 (2009).
60. M. S. Klausen, *et al.*, NetSurfP-2.0: Improved prediction of protein structural features by integrated deep learning. *Proteins: Structure, Function, and Bioinformatics* **87**, 520–527 (2019).
61. N. V. Chawla, K. W. Bowyer, L. O. Hall, W. P. Kegelmeyer, SMOTE: Synthetic Minority Over-sampling Technique. *Journal of Artificial Intelligence Research* **16**, 321–357 (2002).
62. S. Wang, *et al.*, Peptides with selective affinity for carbon nanotubes. *Nature Mater* **2**, 196–200 (2003).
63. J. Liu, L. Yang, A. J. Hopfinger, Affinity of Drugs and Small Biologically Active Molecules to Carbon Nanotubes: A Pharmacodynamics and Nanotoxicity Factor? *Mol. Pharmaceutics* **6**, 873–882 (2009).
64. M. Vihinen, E. Torkkila, P. Riikonen, Accuracy of protein flexibility predictions. *Proteins: Structure, Function, and Bioinformatics* **19**, 141–149 (1994).
65. M. Saeedimazine, E. G. Brandt, A. P. Lyubartsev, Atomistic Perspective on Biomolecular Adsorption on Functionalized Carbon Nanomaterials under Ambient Conditions. *J. Phys. Chem. B* **125**, 416–430 (2021).
66. H. Chaudhary, *et al.*, Intrinsically disordered protein as carbon nanotube dispersant: How dynamic interactions lead to excellent colloidal stability. *Journal of Colloid and Interface Science* **556**, 172–179 (2019).

67. A. Hirano, T. Kameda, Aromaphilicity Index of Amino Acids: Molecular Dynamics Simulations of the Protein Binding Affinity for Carbon Nanomaterials. *ACS Appl. Nano Mater.* **4**, 2486–2495 (2021).
68. Z. He, J. Zhou, Probing carbon nanotube–amino acid interactions in aqueous solution with molecular dynamics simulations. *Carbon* **78**, 500–509 (2014).
69. V. Zorbas, *et al.*, Importance of Aromatic Content for Peptide/Single-Walled Carbon Nanotube Interactions. *J. Am. Chem. Soc.* **127**, 12323–12328 (2005).
70. S. M. Tomásio, T. R. Walsh, Modeling the Binding Affinity of Peptides for Graphitic Surfaces. Influences of Aromatic Content and Interfacial Shape. *J. Phys. Chem. C* **113**, 8778–8785 (2009).
71. A. A. Alizadehmojarad, *et al.*, Binding Affinity and Conformational Preferences Influence Kinetic Stability of Short Oligonucleotides on Carbon Nanotubes. *Advanced Materials Interfaces* **7**, 2000353 (2020).
72. J. F. Campbell, I. Tessmer, H. H. Thorp, D. A. Erie, Atomic Force Microscopy Studies of DNA-Wrapped Carbon Nanotube Structure and Binding to Quantum Dots. *J. Am. Chem. Soc.* **130**, 10648–10655 (2008).
73. A. G. Beyene, *et al.*, Ultralarge Modulation of Fluorescence by Neuromodulators in Carbon Nanotubes Functionalized with Self-Assembled Oligonucleotide Rings. *Nano Lett.* **18**, 6995–7003 (2018).
74. D. Roxbury, J. Mittal, A. Jagota, Molecular-Basis of Single-Walled Carbon Nanotube Recognition by Single-Stranded DNA. *Nano Lett.* **12**, 1464–1469 (2012).
75. E. S. Jeng, A. E. Moll, A. C. Roy, J. B. Gastala, M. S. Strano, Detection of DNA Hybridization Using the Near-Infrared Band-Gap Fluorescence of Single-Walled Carbon Nanotubes. *Nano Lett.* **6**, 371–375 (2006).
76. F. Schöppler, *et al.*, Molar Extinction Coefficient of Single-Wall Carbon Nanotubes. *J. Phys. Chem. C* **115**, 14682–14686 (2011).
77. F. K. Brunecker, F. Schöppler, T. Hertel, Interaction of Polymers with Single-Wall Carbon Nanotubes. *J. Phys. Chem. C* **120**, 10094–10103 (2016).
78. T. A. Davis, L. A. Holland, Peptide Probe for Multiwalled Carbon Nanotubes: Electrophoretic Assessment of the Binding Interface and Evaluation of Surface Functionalization. *ACS Appl. Mater. Interfaces* **10**, 11311–11318 (2018).
79. D. Nepal, K. E. Geckeler, pH-Sensitive Dispersion and Debundling of Single-Walled Carbon Nanotubes: Lysozyme as a Tool. *Small* **2**, 406–412 (2006).
80. D. W. Horn, K. Tracy, C. J. Easley, V. A. Davis, Lysozyme Dispersed Single-Walled Carbon Nanotubes: Interaction and Activity. *J. Phys. Chem. C* **116**, 10341–10348 (2012).
81. B. D. Holt, M. C. McCorry, P. D. Boyer, K. N. Dahl, M. F. Islam, Not all protein-mediated single-wall carbon nanotube dispersions are equally bioactive. *Nanoscale* **4**, 7425–7434 (2012).
82. K. Matsuura, *et al.*, Selectivity of water-soluble proteins in single-walled carbon nanotube dispersions. *Chemical Physics Letters* **429**, 497–502 (2006).
83. M. Zöller, CD44: can a cancer-initiating cell profit from an abundantly expressed molecule? *Nature Reviews Cancer* **11**, 254–267 (2011).
84. A. Rives, *et al.*, Biological structure and function emerge from scaling unsupervised learning to 250 million protein sequences. *PNAS* **118** (2021).
85. F. Pedregosa, *et al.*, Scikit-learn: Machine Learning in Python. *The Journal of Machine Learning Research* **12**, 2825–2830 (2011).
86. N. Otsu, A Threshold Selection Method from Gray-Level Histograms. *IEEE Transactions on Systems, Man, and Cybernetics* **9**, 62–66 (1979).
87. , XGBoost | Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (April 12, 2021).

88. D. Roxbury, P. V. Jena, Y. Shamay, C. P. Horoszkó, D. A. Heller, Cell Membrane Proteins Modulate the Carbon Nanotube Optical Bandgap via Surface Charge Accumulation. *ACS Nano* **10**, 499–506 (2016).