

1 **CIRFESS: An interactive resource for querying the set of theoretically detectable**
2 **peptides for cell surface and extracellular enrichment proteomic studies**

3

4 Matthew Waas (<https://orcid.org/0000-0003-4537-1502>)¹, Jack Littrell ([https://orcid.org/0000-](https://orcid.org/0000-0003-1264-894X)
5 [0003-1264-894X](https://orcid.org/0000-0003-1264-894X))¹, Rebekah L. Gundry (<https://orcid.org/0000-0002-9263-833X>)¹

6 ¹CardiOmics Program, Center for Heart and Vascular Research; Division of Cardiovascular
7 Medicine; and Department of Cellular and Integrative Physiology, University of Nebraska
8 Medical Center, Omaha, NE, 68198, USA

9

10 **Running Title:** CIRFESS informs cell surface proteomics

11

12

13

14

15

16

17

18 Address reprint requests to: Rebekah L. Gundry, PhD, University of Nebraska Medical Center,
19 Department of Cellular and Integrative Physiology, 985850 Nebraska Medical Center, Omaha,
20 NE, 68198-5850, USA, Telephone: 402-559-4426, Fax: 402-559-4438, Email:
21 rebekah.gundry@unmc.edu

1 **Abstract**

2 Cell surface transmembrane, extracellular, and secreted proteins are high value targets
3 for immunophenotyping, drug development, and studies related to intercellular communication
4 in health and disease. As the number of specific and validated affinity reagents that target this
5 subproteome are limited, mass spectrometry (MS)-based approaches will continue to play a
6 critical role in enabling discovery and quantitation of these molecules. Given the technical
7 considerations that make MS-based cell surface proteome studies uniquely challenging, it can
8 be difficult to select an appropriate experimental approach. To this end, we have integrated
9 multiple prediction strategies and annotations into a single online resource, Compiled Interactive
10 Resource for Extracellular and Surface Studies (CIRFESS). CIRFESS enables rapid
11 interrogation of the human proteome to reveal the cell surface proteome theoretically detectable
12 by current approaches and highlights where current prediction strategies provide concordant
13 and discordant information. We applied CIRFESS to identify the percentage of various subsets
14 of the proteome which are expected to be captured by targeted enrichment strategies, including
15 two established methods and one that is possible but not yet demonstrated. These results will
16 inform the selection of available proteomic strategies and development of new strategies to
17 enhance coverage of the cell surface and extracellular proteome. CIRFESS is available at
18 www.cellsurfer.net/cirfess.

19

1 Introduction

2 The emergence of proteomics as a major discipline within the life science has been in no
3 small part due to the development and eager adoption of computational strategies to enable the
4 rapid analysis of mass spectrometry (MS) data files and inferred biological results. Since 1994,
5 when the Yates laboratory introduced SEQUEST¹, the first computational tool for fully
6 automated database searching, continued developments in database construction, algorithm
7 design, and software development have propelled the evolution of MS-based proteomics²⁻⁹. All
8 aspects of MS-based proteomics, including interpretation of raw spectra and database
9 searching, visualization of results, and subsequent biological inferences benefit from advances
10 in bioinformatics. Beyond the analysis of experimental data, data science tools that integrate
11 machine-learning or ontological resources have become increasingly popular for prediction and
12 classification of protein-level information, a subject of recent review¹⁰. Such approaches rely on
13 experimental data to train or inform predictions and annotations, and in turn, the prediction
14 strategies can benefit experimental design.

15 To scientists at the bench, perhaps the most exciting and impactful bioinformatic tools
16 are those that can inform the next experiment. To this end, web-based formats have become
17 increasingly popular resources as they often require less setup (*e.g.* installation), avoid
18 operating system compatibility issues, and can be used in a familiar framework. Hundreds of
19 web-based bioinformatics tools are now available for proteomics (*e.g.*
20 www.expasy.org/proteomics). Current tools span a broad range of utility, including systems-
21 level distribution of proteins based on experimental observations, visualization of experimental
22 results, and prediction and cataloging of specific post-translational modifications, interactomes,
23 and subcellular proteomes. Despite the increase in availability of web-based proteomics tools,
24 there are currently relatively few resources designed to specifically assist in experimental design
25 and analysis of the cell surface proteome.

1 The cell surface and extracellular space contain proteins which play key roles in a wide
2 range of biological processes and can be utilized as valuable markers for immunophenotyping
3 and drug targets. Despite their importance, the cell surface proteome remains relatively poorly
4 characterized compared to the depth that most intracellular proteomes have been described.
5 Given the relative low abundance, presence of hydrophobic transmembrane spanning regions,
6 and dynamic nature of the cell surface proteome owing to continuous cycling of proteins due to
7 internalization, secretion, and stimulus-triggered recruitment to the plasma membrane,
8 specialized techniques are typically required to enhance the detection of cell surface proteins by
9 MS. Such proteomic approaches include enrichment strategies which exploit the biophysical
10 properties of membranes - such as density gradient flotation, differential centrifugation, or silica-
11 bead capture¹¹⁻¹⁵ - or affinity-based approaches that use proximity labels¹⁶⁻¹⁸, lectins¹⁹,
12 metabolic²⁰⁻²² or chemical labels²³⁻²⁷ to enrich cell surface proteins. Application of these
13 approaches have supported efforts to catalog the cell surface and secretome and have led to
14 large scale efforts in experimentation²⁸ and collation²⁹. As with most proteomic methods,
15 currently available strategies to probe the cell surface and secreted proteome are biased
16 towards proteins that contain specific features (e.g. presence of an *N*-glycosylation site or lysine
17 within an extracellular region of the protein that will generate a detectable peptide after trypsin
18 digestion). Also, the implementation of these approaches can be inconsistent among users,
19 resulting in variability in the specificity (*i.e.* cell surface versus non-cell surface) of enrichment
20 observed. Hence, the integration of bioinformatic predictions with experimental data provides
21 orthogonal means to interpret and filter results. Relevant predictions for surface and
22 extracellular proteins include the presence of signal peptides³⁰⁻³⁶ and transmembrane
23 domains^{30,31,37-39}. Other approaches have applied manual curation, ontological annotations, or
24 machine learning approaches to predict the subset of proteins that are localized to the cell
25 surface and extracellular regions⁴⁰⁻⁴³. However, not all cell surface proteins contain canonical
26 signal peptides⁴⁴. Also, GPI-anchored and extracellular matrix or secreted proteins do not

1 contain transmembrane domains, and gene ontology annotations may be insufficiently specific
2 (e.g. cell surface versus membrane). Thus, filtering a proteomic dataset by these constraints
3 often does not provide the complete picture of the cell surface proteome for a specific cell type.

4 Based on these limitations, in theory, it remains necessary to rely on experimental data
5 to precisely define the proteins localized to the cell surface in a specific cell type. This leads to
6 the question: *Which experimental approach is the best to use?* No doubt, the answer will be
7 context dependent. If a specific monoclonal antibody is available, flow cytometry can be an
8 effective approach for determining the surface localization of a protein. If antibodies are not
9 available, the MS-based proteomic method of choice will depend on whether a cell surface
10 proteome-wide screen or detection of a particular protein or subclass of proteins is desired. It
11 will depend on the type and availability of the source material (e.g. metabolic labeling
12 approaches cannot be used routinely for the analysis of primary human cells). Given the
13 numerous technical considerations that make cell surface proteome studies uniquely
14 challenging, it can be difficult to decide which approach to use. Currently, there is no single
15 bioinformatic tool that can assist the investigator in determining which MS-based method is
16 likely to be the most suitable approach for surface proteome studies.

17 To address this, we constructed a resource that integrates multiple prediction strategies
18 and annotations relevant for the analysis of cell surface and extracellular proteins by MS and
19 applied it to interrogate the human proteome. The results from these resources were compiled
20 into a single interface and are accessible via a web-application termed Compiled Interactive
21 Resource for Extracellular and Surface Studies (CIRFESS), accessible at
22 www.cellsurfer.net/cirfess. By bringing together key resources used to interrogate the surface
23 and extracellular space, CIRFESS helps to prevent duplication of efforts and continued ping-
24 pong of separate prediction servers for the same set of proteins. We expect CIRFESS will be

1 informative for a broad range of future applications and will inform the selection or development
 2 of proteomic strategies to enhance coverage of the cell surface and extracellular proteome.

3 **Methods**

4 *Database and prediction server access*

5 The human reference proteome was downloaded from UniProt (canonical only, 20416 entries,
 6 accessed Sept 13, 2019). The proteome was filtered and split to meet the requirements of the
 7 individual prediction servers (e.g. length of proteins, number of entries). Default scoring settings
 8 were applied, and the outputs were collected as specified: TMHMM– ‘one line per protein’,
 9 Phobius – ‘Short’, PrediSi, – ‘Text’, Signal P – ‘Short output’. For analysis involving the protein-
 10 level evidence, the “Protein existence” field for each accession number was retrieved from
 11 UniProt. To aid in interpretation, a summary of the different categories is provided in the Table 1
 12 (adapted from https://www.uniprot.org/help/protein_existence).

Assigned Bin	UniProt “Protein existence” value(s)	UniProt Definition
Protein-level evidence	Experimental evidence at protein level	There is clear experimental evidence for the existence of the protein. The criteria include partial or complete Edman sequencing, clear identification by mass spectrometry, X-ray or NMR structure, good quality protein-protein interaction or detection of the protein by antibodies.
Transcript-level evidence	Experimental evidence at transcript level	The existence of a protein has not been strictly proven but that expression data (such as existence of cDNA(s), RT-PCR or Northern blots) indicate the existence of a transcript.
Other	Protein inferred by homology	The existence of a protein is probable because clear orthologs exist in closely related species.
	Protein predicted	Entries without evidence at protein, transcript, or homology levels.
	Protein uncertain	The existence of the protein is unsure.

13 **Table 1:** Summary of UniProt levels of “Protein existence” and the corresponding bins used in
 14 the analysis of levels of evidence for the different subsets of proteins.

15

16 *Integrating prediction outputs*

17 The outputs from the independent prediction servers were parsed and integrated using Python
 18 3.7. The source code is made available as a Jupyter Notebook to enable implementation on

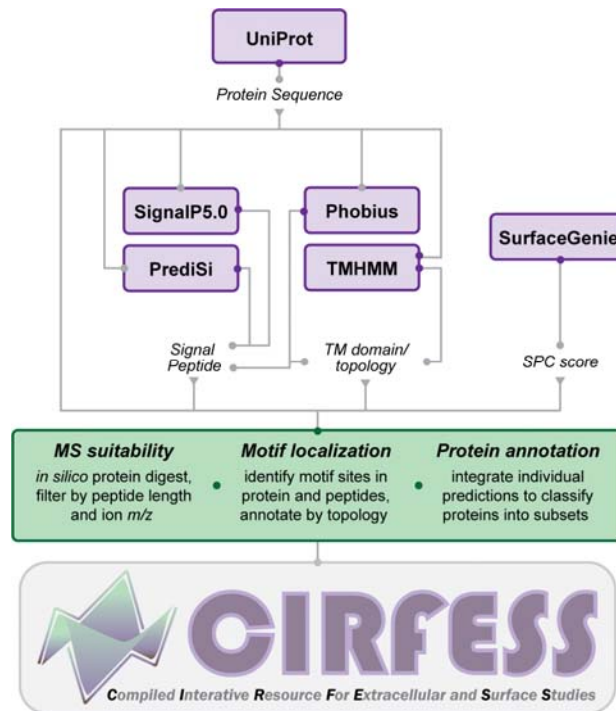
1 batch predictions from TMHMM, Phobius, PrediSi and SignalP^{30,36,39,45} for other species with
2 minimal alteration (<https://github.com/GundryLab/cirfess>). A schematic of the inputs and the
3 generated data structure is shown in Figure 1.

4 *Evaluating peptides, motifs, and topology*

5 Protein sequences were digested *in silico* to generate a list of potential peptides using the
6 canonical tryptic cleavage site, X[R/K] where X is not P, allowing for up to two missed
7 cleavages. This list of peptides is subsequently annotated with the following information: (1)
8 presence of motifs for relevant proteomic capture strategies; N[!P][S/T/C/V] for *N*-glycan based
9 capture, C for cysteine-based capture, K for lysine-based capture, (2) topological information –
10 which residues and motifs are predicted to be intracellular and extracellular, and (3) suitability
11 for a standard bottom-up proteomic experiment – length > 5, *m/z* of 2+ or 3+ charge state
12 peptide < 2000.

13 *CIRFESS Web application*

14 A web application (CIRFESS) for accessing the data structure containing the parsed prediction
15 outputs was developed in R⁴⁶ using the Shiny package and is available at
16 www.cellsurfer.net/cirfess. Source code is available at <https://github.com/GundryLab/cirfess>.



1
2 **Figure 1:** Schematic of the resources used (in purple) and the analyses performed (in green) for
3 the construction of CIRFESS.

4

5 *Statistical Analysis*

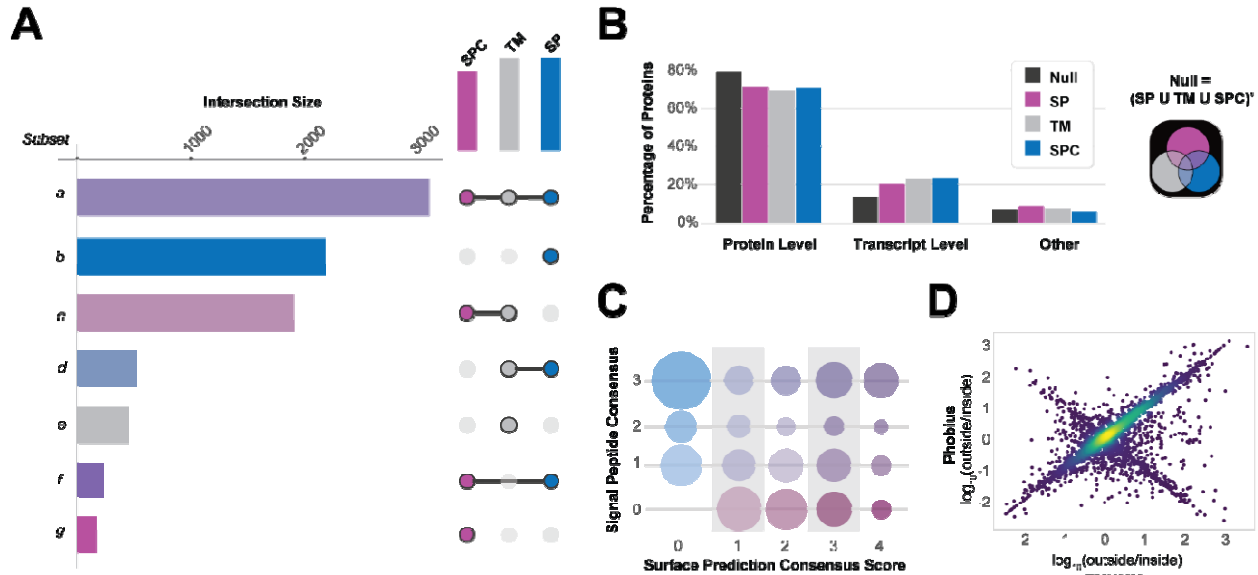
6 Chi-square analyses were performed using the *chisq.test* function in R (version 3.6.2). Students
7 t-test were performed using the *t.test* function in Excel.

8 **Results and Discussion**

9 *Proteome-wide comparison of prediction strategies*

10 There are multiple sources of information to consider for evidence of surface or extracellular
11 localization. Here we utilized transmembrane (TM) predictions^{30,36}, signal peptide (SP)
12 predictions^{31,39}, and the Surface Prediction Consensus (SPC)⁴⁷ score. To highlight the
13 complementarity of these measures and justification for inclusion in this resource, we performed
14 set analysis of the human proteome using the UpSetR⁴⁸ web application. As shown in Figure
15 2A, the combination of these measures reveals distinct subsets of proteins. For instance,

1 proteins which have a SP, but no TM domain (subsets b,f) are often considered to be the set of
 2 secreted proteins. For the set of proteins with predicted TM domains and SP, the SPC score is
 3 helpful for distinguishing between organelle membrane proteins (subset d) and cell surface
 4 proteins (subset a). Finally, proteins with an SPC score but not TM or SP may contain GPI-



5 anchored proteins or those without a canonical SP (subset g).

6

7 **Figure 2:** Contrasting views of the human proteome based on prediction strategies relevant for
 8 the cell surface proteome. (A) UpSet plot illustrating overlap in human proteins that are
 9 classified as containing a SP, TM domain, or SPC > 0. (B) Bar graph depicting the percentage
 10 of the human proteome with different levels of evidence, gathered from the "Protein existence"
 11 level listed for each accession number in UniProt, that are classified as containing a SP, TM
 12 domain, or SPC >0, or none of these features (Null). (C) Relationship between different levels of
 13 consensus for SP prediction and SPC score. Here, SP prediction consensus was calculated in a
 14 manner analogous to SPC score, where the number of positive SP predictions from SignalP,
 15 PrediSI, or Phobius was summed to generate a consensus score ranging from 0 to 3. (D) Plot
 16 depicting the \log_{10} ratio of extracellular to intracellular residues predicted by TMHMM and
 17 Phobius highlighting that the opposite orientation is predicted for a subset of proteins.

18

19 To assess the level of experimental data that currently exists for cell surface proteins, we
 20 considered the "Protein existence" annotation within UniProt. For these and further analyses, we

1 compared proteins with positive SP predictions, positive TM predictions, or those with SPC
 2 scores > 0, with proteins that were negative for all three analysis – which we term the ‘Null set’
 3 of proteins (shown visually in Figure 2B). Compared to the Null set, all three classes of proteins
 4 have a lower percentage of members with protein-level evidence, 79% for Null proteins
 5 compared to 71%, 69%, and 71% for SP, TM, and SPC proteins, respectively (Figure 2B). The
 6 difference between the observed frequencies of the different levels of evidence were significant
 7 between Null and each other subset of proteins (SP, TM, and SPC), as revealed by Chi-square
 8 testing (with p-values of $<2.2 \times 10^{-16}$ for each test). Though mass spectrometry is not the only
 9 source of protein-level evidence for UniProt (see Methods), a potential explanation for this
 10 discrepancy is a statistical difference in the number of MS-suitable peptides between these
 11 subsets, as revealed by Student’s t-test summarized in Table 2. Nevertheless, this analysis
 12 highlights the need for further experimental investigation of the cell surface proteome as this
 13 class is less well-represented by experimental evidence than other subproteomes.

Missed Cleavages	Null			SP			TM			SPC		
	0	≤ 1	≤ 2	0	≤ 1	≤ 2	0	≤ 1	≤ 2	0	≤ 1	≤ 2
Mean # of Peptides / protein	33.1	87.5	146.4	28.3	70.5	112.3	26.2	65.0	103.4	28.1	69.8	111.0
Median # of Peptides / protein	25	66	109	20	50	79	19	48	76	20	50	79
t-test p-value (compared to Null-set)	-	-	-	6×10^{-23}	1×10^{-41}	2×10^{-58}	2×10^{-47}	8×10^{-75}	1×10^{-97}	2×10^{-23}	4×10^{-42}	4×10^{-59}

14 **Table 2:** Summary of the average number of peptides per protein that are “ok for MS” in
 15 different subsets of proteins. The t-test p-values were calculated by comparing the distribution to
 16 the Null-set of numbers of peptides (with the corresponding amount of numbers of missed
 17 cleavages).

18
 19 Another benefit of integrating these disparate predictions into a single analysis is
 20 revealed by looking at examples for which they do and do not agree. For example, stratifying
 21 proteins with positive SP predictions (6026 proteins) by the number of algorithms for which it
 22 was positive reveals that slightly over half (3192, 53%) are predicted by all 3 algorithms. Here,

1 SP prediction consensus was calculated in manner analogous to SPC score, where the number
2 of positive SP predictions from SignalP, PrediSI, or Phobius was summed to generate a
3 consensus score ranging from 0 to 3. Plotting the number of positive SP predictions against
4 SPC score reveals a positive relationship between SPC score and number of positive SP
5 predictions for proteins with SPC score >0 (Figure 2C). However, it also reveals that the majority
6 of proteins with three positive SP predictions has an SPC score of 0 (1648 of 3192, 51.6%).
7 This suggests that secreted proteins may contain signal peptide sequences that are easier to
8 recognize by prediction algorithms than proteins translocated through the membrane.

9 Focusing on TM proteins, TMHMM and Phobius predict 5353 and 5471 proteins with TM
10 domains, respectively, with 4846 proteins in common and 1132 proteins unique to a single
11 prediction strategy (507 and 625 in TMHMM and Phobius, respectively). While overall there is
12 strong consensus between the two algorithms for predicting which proteins contain TM
13 domains, the number of TM domains predicted differs for 1306 out of the 4846 commonly
14 predicted proteins. Furthermore, the opposite membrane orientation was predicted for a subset
15 of proteins, visualized by plotting the \log_{10} ratio of the predicted extracellular to intracellular
16 residues (Figure 2D). Altogether, these analyses demonstrate the value of integrating data from
17 multiple sources and reveal that no single feature is sufficient to comprehensively predict the set
18 of cell surface and extracellular proteins.

19 *Motif coverage of extracellular and surface predicted proteins*

20 As prediction strategies alone are insufficient to define the set of proteins localized to the cell
21 surface and extracellular space, experimentation is required. To aid in the selection of proteomic
22 strategies that are likely to produce the desired coverage of the cell surface proteome, the SP,
23 TM, and SPC analyses described above were integrated with *in silico* analyses designed to
24 predict which proteins would generate tryptic peptides likely to be detectable by electrospray
25 MS, and of those, which are expected to be captured by application of commonly used

1 biorthogonal enrichment strategies targeting *N*-glycans and lysines^{23–25,27,49–54}. We also
2 considered cysteines as they are enriched in surface proteins compared to nonsurface
3 proteins⁴⁰ and numerous affinity reagents are available for targeting cysteines⁵⁵, although this is
4 not yet a widely described approach for cell surface proteins. Important for the *N*-glycan
5 approach, although strategies that specifically enrich peptides from the extracellular space
6 (glycan biotinylation is performed on cells with intact plasma membranes) provide an additional
7 level of experimental evidence for surface localization, it is possible to capture *N*-glycopeptides
8 from whole cell lysate. In this case, the MS-based evidence for a glycan modifying an
9 asparagine within the consensus motif for *N*-glycosylation is proposed to serve as standalone
10 evidence for surface localization. Canonically, the consensus motif has been described as
11 NXS/T where X is any amino acid except proline. However, more recently, evidence for *N*-
12 glycosylation has been put forth at NXC^{56,57} and NXV⁵⁸. Here, we investigated the frequency of
13 the various consensus motifs occurring in SP, TM, SPC and Null (meaning the protein contains
14 no SP, TM, or SPC) sets of proteins. First, the probability of each motif occurring within the
15 subset of proteins was calculated with respect to the amino acid frequencies. The expected
16 frequencies based on amino acid compositions was consistent among the sets of proteins for
17 each motif (0.26 ± 0.003 %, 0.18 ± 0.009 %, 0.09 ± 0.009 %, and 0.21 ± 0.019 % for NXS, NXT,
18 NXC, and NXV respectively). Next, the observed frequency of each motif was calculated for
19 each subset of proteins. The natural log of the odds ratio of observed to expected for each
20 subset of proteins was calculated and plotted for each motif as well as the canonical (NXS/T)
21 and complete consensus motifs (NXS/T/C/V) (Figure 3A). The results reveal that the NXS and
22 NXT occur more frequently than expected and NXC and NXV occur less frequently than
23 expected for SP, TM, and SPC proteins. Whereas NXS and NXT occur at about the expected
24 rate for the Null set of proteins, NXC and NXV occur slightly above the expected frequency.
25 While the complete consensus motif occurs more frequently than expected for SP, TM, and
26 SPC proteins, the canonical motif demonstrates a much higher odds ratio, especially relative to

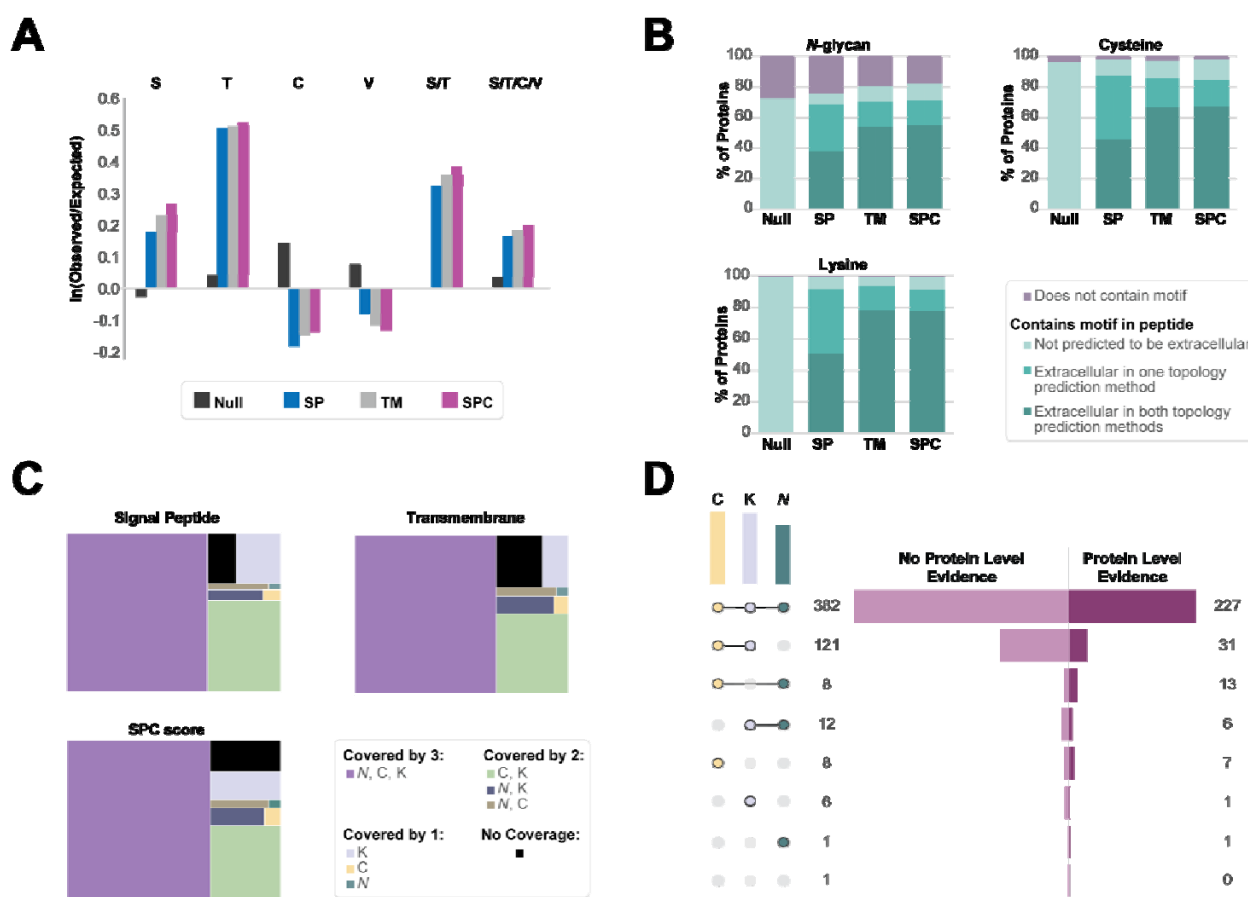
1 the odds-ratio for the Null set. This analysis suggests that while the mere presence of the
2 consensus motif provides some evidentiary weight to the localization of a surface protein – (1) it
3 should not be considered conclusive, and (2) the canonical motif provides more meaningful
4 information than the complete consensus motif. Based on these analyses, we elected to only
5 consider the canonical motif as potential targets for *N*-glycan capture for subsequent analyses.

6

7 By integrating the topology information provided by the TM predictions with the locations of
8 motifs within proteins, we estimated the coverage that each capture strategy would provide for
9 each subset of the proteome (Null, SP, TM and SPC). For this analysis, proteins were
10 categorized based on whether they contained the relevant motif and whether the motif was in a
11 region determined to be extracellular by one or both TM prediction strategies. The percentage of
12 proteins for which a predicted extracellular motif was located within an MS-suitable peptide was
13 recorded (Figure 3B). This analysis revealed that while 72% of Null proteins contain a
14 consensus motif for *N*-glycosylation, none of the glycopeptides are predicted to be in the
15 extracellular domain. In contrast, of the SPC proteins which contain the consensus motif, 86%
16 of those proteins contain at least one peptide contains the consensus motif within the predicted
17 extracellular domain. These results were further summarized by calculating the percentage of
18 each subproteome that is predicted to be covered by each or multiple capture strategies (Figure
19 3C). Overall, querying the results from this analysis provides a strategy for investigators to
20 rapidly interrogate the human proteome to determine which experimental strategy is most likely
21 to be useful to address their biological question. In summary, $66.4 \pm 0.4\%$ of SP, TM and SPC
22 proteins are likely to be captured by any of the three strategies, $17.3 \pm 1.8\%$ are detectable by
23 cysteine or lysine capture, but not detectable by *N*-glycan strategies, $7.4 \pm 1.4\%$ are detectable
24 by a single strategy, and $5.7 \pm 1.2\%$ are not detectable by any of the three strategies considered
25 here. The identity of the proteins within each classification are provided in Supplemental Table

1 1A-C and these results provide actionable data related to high interest targets. For example, of
 2 the 825 human G-protein coupled receptors (GPCR), a striking 65.3% lack protein-level
 3 evidence within UniProt. Of these, all but one are predicted to be captured by at least one
 4 enrichment strategy and 70.9% of them are predicted to be captured by all three strategies.
 5 Supplementary Table 1D contains the identity of the GPCR proteins and which enrichment
 6 strategies are predicted to capture them.

7



8

9 **Figure 3.** Results of CIRFESS analysis of the human proteome to assess predicted coverage
 10 provided by three common cell surface proteomic enrichment strategies. (A) The natural log of
 11 the odds ratio for observed-to-expected frequency of each permutation of the *N*-glycan
 12 consensus motif along with the canonical (S/T) and complete (S/T/C/V) consensus motif. (B)
 13 The expected coverage of the different subsets of proteins for each enrichment strategy broken
 14 down by which proteins have peptides with predicted extracellular motifs by one or both TM
 15 prediction methods. (C) The makeup of SP, TM, and SPC score proteins based on the

1 overlapping coverage of the three individual enrichment strategies. (D) The set of human
2 GPCRs based on expected coverage for enrichment strategy and level of evidence in UniProt.

3

4

5

6 *Critical Considerations*

7 Results from the current implementation of CIRFESS are limited to human proteins digested
8 with trypsin and the resulting peptides are detectable in the 2+ or 3+ charge state. These criteria
9 were selected based on common implementation of bottom-up proteomic methods. However, all
10 source files and code are publicly available in the Github repository and a user-specific version
11 of CIRFESS could be generated, requiring minimum alteration to change the *in silico* digestion
12 strategies or criteria for MS-compatible peptide filtering. Implementation on other species would
13 require submission of proteins to the individual prediction servers, but the source code includes
14 scripts to parse and integrate the generated output files. Another critical assumption is related to
15 the *N*-glycan capture strategy where detection depends on the glycosite being occupied by a
16 glycan which is sensitive to the oxidation strategy applied (*e.g.* cis diols for meta-periodate⁵⁹).
17 Currently, as it is not possible to predict which sites will be occupied with specific glycan
18 structures, the peptides predicted to be observable by this strategy should be considered with
19 this caveat in mind. It is possible that post-translational modifications may interfere with the
20 digestion, capture, ionization, and identification of peptides in any of the strategies, and
21 therefore experimental observations may not be fully predictable by this bioinformatics
22 approach. Among the post-translation modification which may interfere with cysteine-based
23 capture are disulfide bridges, which were ignored in this analysis, but a reduction step could be
24 included prior to labeling in such an approach. Moreover, for enrichment strategies which use
25 cleavable linkers, residual portions of the linker that remain after cleavage will increase the

1 mass of the resulting peptide. However, the 2000 *m/z* range used here for predicting detectable
2 peptides should accommodate most commonly used reagents. Finally, it may be beneficial to
3 combine the results from CIRFESS analysis with predictions for peptide detectability⁶⁰⁻⁶² or
4 proteotypicity⁶³ to better inform the set of peptides which are most likely to be observable or
5 informative.

6 **Conclusion**

7 CIRFESS is a web-based tool designed to accelerate cell surface proteome studies by
8 eliminating the need to query each bioinformatics source separately and integrating disparate
9 features into a single streamlined resource and output. Within the CIRFESS interface, users
10 are able to perform single and batch querying of protein accession numbers to extract protein-
11 level and peptide-level annotations as well as information about numbers of motifs and motif-
12 containing peptides. Results may be queried for proteins or protein classes of interest to inform
13 the design of the next experiment. We anticipate that CIRFESS will be broadly applicable for
14 multiple applications across a broad range of biology and disease studies. While there still exist
15 significant technical challenges associated with the implementation of these technologies,
16 particularly on sample-limited systems, these analyses suggest that acquiring protein-level
17 evidence for the majority of predicted cell surface proteins is a matter of applying the right
18 technology to a relevant biological system. Overall, we expect CIRFESS will promote the
19 rational selection of the most apt cell surface proteomic methods and will inspire continued
20 method development (e.g. cysteine-targeting) to enable detection of the human proteome not
21 predicted to be accessible by established surface protein enrichment methods.

22 **Author Contributions**

23 R.L.G. and M.W. conceived the study; R.L.G. supervised the study; M.W. and J. L. developed
24 the algorithms and developed the web application; M.W. and R.L.G. analyzed data; M.W.

1 generated figures; M.W. and R.L.G. co-wrote the manuscript; All authors approved the final
2 manuscript.

3 **Acknowledgements**

4 This work was supported by the National Institutes of Health [R01-HL126785, R01-HL134010,
5 P20GM104320 (as a pilot grant award) to R.L.G.; F31-HL140914 to M.W.]; and JDRF [2-SRA-
6 2019-829-S-B to R.L.G.]; Special thanks to Dr. Christopher Ashwood and Linda Berg Luecke for
7 critical review of the manuscript and insightful discussions and Dr. David Tabb for helpful
8 consultation. Funding sources were not involved in study design, data collection, interpretation,
9 analysis or publication.

10

11 **Supplementary Information**

12 **Supplemental Table 1.** A. Coverage of SP proteins. B. Coverage of TM proteins. C. Coverage
13 of SPC Proteins. D. Coverage and evidence of G-protein coupled receptors.

14

1 References

- 2 (1) Eng, J. K.; McCormack, A. L.; Yates, J. R. An Approach to Correlate Tandem Mass
3 Spectral Data of Peptides with Amino Acid Sequences in a Protein Database. *Journal of*
4 *the American Society for Mass Spectrometry* **1994**, *5* (11), 976–989.
- 5 (2) Eng, J. K.; Jahan, T. A.; Hoopmann, M. R. Comet: An Open-Source MS/MS Sequence
6 Database Search Tool. *Proteomics* **2013**, *13* (1), 22–24.
- 7 (3) Gluck, F.; Hoogland, C.; Antinori, P.; Robin, X.; Nikitin, F.; Zufferey, A.; Pasquarello, C.;
8 Fétaud, V.; Dayon, L.; Müller, M.; Lisacek, F.; Geiser, L.; Hochstrasser, D.; Sanchez, J.
9 C.; Scherl, A. EasyProt - An Easy-to-Use Graphical Platform for Proteomics Data
10 Analysis. *Journal of Proteomics* **2013**, *79*, 146–160.
- 11 (4) Wenger, C. D.; Coon, J. J. A Proteomics Search Algorithm Specifically Designed for
12 High-Resolution Tandem Mass Spectra. *Journal of Proteome Research* **2013**, *12* (3),
13 1377–1386.
- 14 (5) Dorfer, V.; Pichler, P.; Stranzl, T.; Stadlmann, J.; Taus, T.; Winkler, S.; Mechtler, K. MS
15 Amanda, a Universal Identification Algorithm Optimized for High Accuracy Tandem Mass
16 Spectra. *Journal of Proteome Research* **2014**, *13* (8), 3679–3684.
- 17 (6) Kim, S.; Pevzner, P. A. MS-GF+ Makes Progress towards a Universal Database Search
18 Tool for Proteomics. *Nature Communications* **2014**, *5*.
- 19 (7) Risk, B. A.; Spitzer, W. J.; Giddings, M. C. Peppy: Proteogenomic Search Software.
20 *Journal of Proteome Research* **2013**, *12* (6), 3019–3025.
- 21 (8) Cox, J.; Mann, M. MaxQuant Enables High Peptide Identification Rates, Individualized
22 p.p.b.-Range Mass Accuracies and Proteome-Wide Protein Quantification. *Nature*
23 *Biotechnology* **2008**, *26* (12), 1367–1372.
- 24 (9) Perkins, D. N.; Pappin, D. J. C.; Creasy, D. M.; Cottrell, J. S. Probability-Based Protein
25 Identification by Searching Sequence Databases Using Mass Spectrometry Data. In
26 *Electrophoresis*; Wiley-VCH Verlag, 1999; Vol. 20, pp 3551–3567.
- 27 (10) Ching, T.; Himmelstein, D. S.; Beaulieu-Jones, B. K.; Kalinin, A. A.; Do, B. T.; Way, G. P.;
28 Ferrero, E.; Agapow, P. M.; Zietz, M.; Hoffman, M. M.; Xie, W.; Rosen, G. L.; Lengerich,
29 B. J.; Israeli, J.; Lanchantin, J.; Woloszynek, S.; Carpenter, A. E.; Shrikumar, A.; Xu, J.;
30 Cofer, E. M.; Lavender, C. A.; Turaga, S. C.; Alexandari, A. M.; Lu, Z.; Harris, D. J.;
31 Decaprio, D.; Qi, Y.; Kundaje, A.; Peng, Y.; Wiley, L. K.; Segler, M. H. S.; Boca, S. M.;
32 Swamidass, S. J.; Huang, A.; Gitter, A.; Greene, C. S. Opportunities and Obstacles for
33 Deep Learning in Biology and Medicine. *Journal of the Royal Society Interface* **2018**, *15*
34 (141).
- 35 (11) Chaney, L. K.; Jacobson, B. S. Coating Cells with Colloidal Silica for High Yield Isolation
36 of Plasma Membrane Sheets and Identification of Transmembrane Proteins. *J. Biol.*
37 *Chem.* **1983**, *258* (16), 10062–10072.
- 38 (12) Kim, Y.; Elschenbroich, S.; Sharma, P.; Sepiashvili, L.; Gramolini, A. O.; Kislinger, T. Use
39 of Colloidal Silica-Beads for the Isolation of Cell-Surface Proteins for Mass Spectrometry-
40 Based Proteomics. In *Immune Receptors: Methods and Protocols*; Rast, J. P., Booth, J.
41 W. D., Eds.; Methods in Molecular Biology; Humana Press: Totowa, NJ, 2011; pp 227–
42 241.
- 43 (13) Mulvey, C. M.; Breckels, L. M.; Geladaki, A.; Britovšek, N. K.; Nightingale, D. J. H.;
44 Christoforou, A.; Elzek, M.; Deery, M. J.; Gatto, L.; Lilley, K. S. Using HyperLOPIT to
45 Perform High-Resolution Mapping of the Spatial Proteome. *Nature Protocols* **2017**, *12*
46 (6), 1110–1135.
- 47 (14) Jett, M.; Seed, T. M.; Jamieson, G. A. Isolation and Characterization of Plasma
48 Membranes and Intact Nuclei from Lymphoid Cells. *J. Biol. Chem.* **1977**, *252* (6), 2134–
49 2142.

- 1 (15) Jones, D. H.; Matus, A. I. Isolation of Synaptic Plasma Membrane from Brain by
2 Combined Flotation-Sedimentation Density Gradient Centrifugation. *Biochimica et*
3 *Biophysica Acta (BBA) - Biomembranes* **1974**, 356 (3), 276–287.
- 4 (16) Rees, J. S.; Li, X.-W.; Perrett, S.; Lilley, K. S.; Jackson, A. P. Selective Proteomic
5 Proximity Labeling Assay Using Tyramide (SPPLAT): A Quantitative Method for the
6 Proteomic Analysis of Localized Membrane-Bound Protein Clusters. *Current Protocols in*
7 *Protein Science* **2015**, 80 (1), 19.27.1-19.27.18.
- 8 (17) Lee, S.-Y.; Kang, M.-G.; Park, J.-S.; Lee, G.; Ting, A. Y.; Rhee, H.-W. APEX
9 Fingerprinting Reveals the Subcellular Localization of Proteins of Interest. *Cell Rep* **2016**,
10 15 (8), 1837–1847.
- 11 (18) Roux, K. J.; Kim, D. I.; Raida, M.; Burke, B. A Promiscuous Biotin Ligase Fusion Protein
12 Identifies Proximal and Interacting Proteins in Mammalian Cells. *J. Cell Biol.* **2012**, 196
13 (6), 801–810.
- 14 (19) Ghosh, D.; Krokhin, O.; Antonovici, M.; Ens, W.; Standing, K. G.; Beavis, R. C.; Wilkins,
15 J. A. Lectin Affinity as an Approach to the Proteomic Analysis of Membrane
16 Glycoproteins. *J. Proteome Res.* **2004**, 3 (4), 841–850.
- 17 (20) Breidenbach, M. A.; Gallagher, J. E. G.; King, D. S.; Smart, B. P.; Wu, P.; Bertozzi, C. R.
18 Targeted Metabolic Labeling of Yeast N-Glycans with Unnatural Sugars. *Proc. Natl. Acad.*
19 *Sci. U.S.A.* **2010**, 107 (9), 3988–3993.
- 20 (21) Smeekens, J. M.; Chen, W.; Wu, R. Mass Spectrometric Analysis of the Cell Surface N-
21 Glycoproteome by Combining Metabolic Labeling and Click Chemistry. *J. Am. Soc. Mass*
22 *Spectrom.* **2015**, 26 (4), 604–614.
- 23 (22) Sun, T.; Yu, S.-H.; Zhao, P.; Meng, L.; Moremen, K. W.; Wells, L.; Steet, R.; Boons, G.-J.
24 One-Step Selective Exoenzymatic Labeling (SEEL) Strategy for the Biotinylation and
25 Identification of Glycoproteins of Living Cells. *J. Am. Chem. Soc.* **2016**, 138 (36), 11575–
26 11582.
- 27 (23) Kalxdorf, M.; Gade, S.; Eberl, H. C.; Bantscheff, M. Monitoring Cell-Surface N-
28 Glycoproteome Dynamics by Quantitative Proteomics Reveals Mechanistic Insights into
29 Macrophage Differentiation. *Molecular & cellular proteomics*: MCP **2017**, 16 (5), 770–
30 785.
- 31 (24) Turtoi, A.; Dumont, B.; Greffe, Y.; Blomme, A.; Mazzucchelli, G.; Delvenne, P.; Mutijima,
32 E. N.; Lifrange, E.; De Pauw, E.; Castronovo, V. Novel Comprehensive Approach for
33 Accessible Biomarker Identification and Absolute Quantification from Precious Human
34 Tissues. *Journal of Proteome Research* **2011**, 10 (7), 3160–3182.
- 35 (25) Wollscheid, B.; Bausch-Fluck, D.; Henderson, C.; O'Brien, R.; Bibel, M.; Schiess, R.;
36 Aebersold, R.; Watts, J. D. Mass-Spectrometric Identification and Relative Quantification
37 of N-Linked Cell Surface Glycoproteins. *Nature biotechnology* **2009**, 27 (4), 378–386.
- 38 (26) Zhang, H.; Li, X. jun; Martin, D. B.; Aebersold, R. Identification and Quantification of N-
39 Linked Glycoproteins Using Hydrazide Chemistry, Stable Isotope Labeling and Mass
40 Spectrometry. *Nature Biotechnology* **2003**, 21 (6), 660–666.
- 41 (27) Castronovo, V.; Kischel, P.; Guillonneau, F.; de Leval, L.; Deféchereux, T.; De Pauw, E.;
42 Neri, D.; Waltregny, D. Identification of Specific Reachable Molecular Targets in Human
43 Breast Cancer Using a Versatile Ex Vivo Proteomic Method. *Proteomics* **2007**, 7 (8),
44 1188–1196.
- 45 (28) Bausch-Fluck, D.; Hofmann, A.; Bock, T.; Frei, A. P.; Cerciello, F.; Jacobs, A.; Moest, H.;
46 Omasits, U.; Gundry, R. L.; Yoon, C.; Schiess, R.; Schmidt, A.; Mirkowska, P.; Härtlová,
47 A.; Van Eyk, J. E.; Bourquin, J.-P.; Aebersold, R.; Boheler, K. R.; Zandstra, P.;
48 Wollscheid, B. A Mass Spectrometric-Derived Cell Surface Protein Atlas. *PLOS ONE*
49 **2015**, 10 (4), e0121314.
- 50 (29) Uhlén, M.; Karlsson, M. J.; Hober, A.; Svensson, A.-S.; Scheffel, J.; Kotol, D.; Zhong, W.;
51 Tebani, A.; Strandberg, L.; Edfors, F.; Sjöstedt, E.; Mulder, J.; Mardinoglu, A.; Berling, A.;

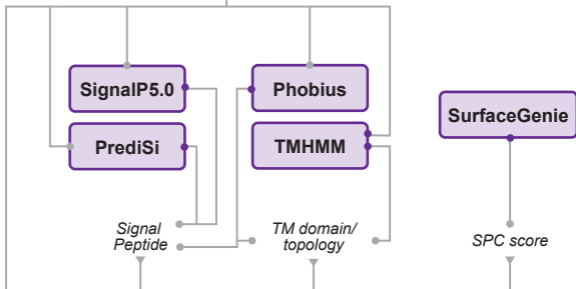
- 1 Ekblad, S.; Dannemeyer, M.; Kanje, S.; Rockberg, J.; Lundqvist, M.; Malm, M.; Volk, A.-
2 L.; Nilsson, P.; Månberg, A.; Dodig-Crnkovic, T.; Pin, E.; Zwahlen, M.; Oksvold, P.; von
3 Feilitzen, K.; Häussler, R. S.; Hong, M.-G.; Lindskog, C.; Ponten, F.; Katona, B.; Vuu, J.;
4 Lindström, E.; Nielsen, J.; Robinson, J.; Ayoglu, B.; Mahdessian, D.; Sullivan, D.; Thul,
5 P.; Danielsson, F.; Stadler, C.; Lundberg, E.; Bergström, G.; Gummesson, A.; Voldborg,
6 B. G.; Tegel, H.; Hober, S.; Forsström, B.; Schwenk, J. M.; Fagerberg, L.; Sivertsson, Å.
7 The Human Secretome. *Science signaling* **2019**, *12* (609).
- 8 (30) Hiller, K.; Grote, A.; Scheer, M.; Münch, R.; Jahn, D. PrediSi: Prediction of Signal
9 Peptides and Their Cleavage Positions. *Nucleic Acids Research* **2004**, *32* (WEB
10 SERVER ISS.).
- 11 (31) Käll, L.; Krogh, A.; Sonnhammer, E. L. L. A Combined Transmembrane Topology and
12 Signal Peptide Prediction Method. *Journal of Molecular Biology* **2004**, *338* (5), 1027–
13 1036.
- 14 (32) Frank, K.; Sippl, M. J. High-Performance Signal Peptide Prediction Based on Sequence
15 Alignment Techniques. *Bioinformatics* **2008**, *24* (19), 2172–2176.
- 16 (33) Signal-3L 3.0: Improving signal peptide prediction through combining attention deep
17 learning with domain rules <http://www.csbio.sjtu.edu.cn/bioinf/Signal-3L/> (accessed Jan
18 19, 2020).
- 19 (34) Reynolds, S. M.; Käll, L.; Riffle, M. E.; Bilmes, J. A.; Noble, W. S. Transmembrane
20 Topology and Signal Peptide Prediction Using Dynamic Bayesian Networks. *PLoS*
21 *Computational Biology* **2008**, *4* (11).
- 22 (35) Savojardo, C.; Martelli, P. L.; Fariselli, P.; Casadio, R. DeepSig: Deep Learning Improves
23 Signal Peptide Detection in Proteins. *Bioinformatics* **2018**, *34* (10), 1690–1696.
- 24 (36) Almagro Armenteros, J. J.; Tsirigos, K. D.; Sønderby, C. K.; Petersen, T. N.; Winther, O.;
25 Brunak, S.; von Heijne, G.; Nielsen, H. SignalP 5.0 Improves Signal Peptide Predictions
26 Using Deep Neural Networks. *Nature Biotechnology* **2019**, *37* (4), 420–423.
- 27 (37) Hirokawa, T.; Boon-Chieng, S.; Mitaku, S. SOSUI: Classification and Secondary Structure
28 Prediction System for Membrane Proteins. *Bioinformatics* **1998**, *14* (4), 378–379.
- 29 (38) Bernhofer, M.; Kloppmann, E.; Reeb, J.; Rost, B. TMSEG: Novel Prediction of
30 Transmembrane Helices. *Proteins: Structure, Function and Bioinformatics* **2016**, *84* (11),
31 1706–1716.
- 32 (39) Sonnhammer, E. L.; von Heijne, G.; Krogh, A. A Hidden Markov Model for Predicting
33 Transmembrane Helices in Protein Sequences. *Proceedings / ... International Conference*
34 *on Intelligent Systems for Molecular Biology*; ISMB. *International Conference on*
35 *Intelligent Systems for Molecular Biology* **1998**, *6*, 175–182.
- 36 (40) Bausch-Fluck, D.; Goldmann, U.; Müller, S.; van Oostrum, M.; Müller, M.; Schubert, O. T.;
37 Wollscheid, B. The in Silico Human Surfaceome. *Proceedings of the National Academy of*
38 *Sciences of the United States of America* **2018**, *115* (46), E10988–E10997.
- 39 (41) da Cunha, J. P. C.; Galante, P. A. F.; de Souza, J. E.; de Souza, R. F.; Carvalho, P. M.;
40 Ohara, D. T.; Moura, R. P.; Oba-Shinja, S. M.; Marie, S. K. N.; Silva, W. A.; Perez, R. O.;
41 Stransky, B.; Pieprzyk, M.; Moore, J.; Caballero, O.; Gama-Rodrigues, J.; Habr-Gama, A.;
42 Kuo, W. P.; Simpson, A. J.; Camargo, A. A.; Old, L. J.; de Souza, S. J. Bioinformatics
43 Construction of the Human Cell Surfaceome. *Proceedings of the National Academy of*
44 *Sciences of the United States of America* **2009**, *106* (39), 16752–16757.
- 45 (42) Town, J.; Pais, H.; Harrison, S.; Stead, L. F.; Bataille, C.; Bunjobpol, W.; Zhang, J.;
46 Rabbitts, T. H. Exploring the Surfaceome of Ewing Sarcoma Identifies a New and Unique
47 Therapeutic Target. *Proceedings of the National Academy of Sciences of the United*
48 *States of America* **2016**, *113* (13), 3603–3608.
- 49 (43) Díaz-Ramos, M. C.; Engel, P.; Bastos, R. Towards a Comprehensive Human Cell-
50 Surface Immunome Database. *Immunology letters* **2011**, *134* (2), 183–187.

- 1 (44) Kuchler, K.; Thorner, J. Membrane Translocation of Proteins without Hydrophobic Signal
2 Peptides. *Current Opinion in Cell Biology* **1990**, 2 (4), 617–624.
- 3 (45) Käll, L.; Krogh, A.; Sonnhammer, E. L. L. A Combined Transmembrane Topology and
4 Signal Peptide Prediction Method. *Journal of Molecular Biology* **2004**, 338 (5), 1027–
5 1036.
- 6 (46) *R Core Team (2014). R: A Language and Environment for Statistical Computing. R*
7 *Foundation for Statistical Computing*; Vienna, Austria.
- 8 (47) Waas, M.; Snarrenberg, S. T.; Littrell, J.; Lipinski, R. A. J.; Hansen, P. A.; Corbett, J. A.;
9 Gundry, R. L. SurfaceGenie: A Web-Based Application for Prioritizing Cell-Type Specific
10 Marker Candidates. <https://doi.org/10.1101/575969> **2020**, In Revision.
- 11 (48) Lex, A.; Gehlenborg, N.; Strobel, H.; Vuillemot, R.; Pfister, H. UpSet: Visualization of
12 Intersecting Sets. *IEEE transactions on visualization and computer graphics* **2014**, 20
13 (12), 1983–1992.
- 14 (49) Gundry, R. L.; Riordon, D. R.; Tarasova, Y.; Chuppa, S.; Bhattacharya, S.; Juhasz, O.;
15 Wiedemeier, O.; Milanovich, S.; Noto, F. K.; Tchernyshyov, I.; Raginski, K.; Bausch-
16 Fluck, D.; Tae, H.-J.; Marshall, S.; Duncan, S. A.; Wollscheid, B.; Wersto, R. P.; Rao, S.;
17 Eyk, J. E. V.; Boheler, K. R. A Cell Surfaceome Map for Immunophenotyping and Sorting
18 Pluripotent Stem Cells. *Molecular & Cellular Proteomics* **2012**, 11 (8), 303–316.
- 19 (50) Ye, X.; Chan, K. C.; Waters, A. M.; Bess, M.; Harned, A.; Wei, B.-R.; Loncarek, J.; Luke,
20 B. T.; Orsburn, B. C.; Hollinger, B. D.; Stephens, R. M.; Bagni, R.; Martinko, A.; Wells, J.
21 A.; Nissley, D. V.; McCormick, F.; Whiteley, G.; Blonder, J. Comparative Proteomics of a
22 Model MCF10A-KRasG12V Cell Line Reveals a Distinct Molecular Signature of the
23 KRasG12V Cell Surface. *Oncotarget* **2016**, 7 (52), 86948–86971.
- 24 (51) Ravenhill, B. J.; Soday, L.; Houghton, J.; Antrobus, R.; Weekes, M. P. Comprehensive
25 Cell Surface Proteomics Defines Markers of Classical, Intermediate and Non-Classical
26 Monocytes. *Scientific Reports* **2020**, 10 (1), 1–11.
- 27 (52) Chauhan, S.; Danielson, S.; Clements, V.; Edwards, N.; Ostrand-Rosenberg, S.;
28 Fenselau, C. Surface Glycoproteins of Exosomes Shed by Myeloid-Derived Suppressor
29 Cells Contribute to Function. *J. Proteome Res.* **2017**, 16 (1), 238–246.
- 30 (53) Chen, W.; Smeekens, J. M.; Wu, R. A Universal Chemical Enrichment Method for
31 Mapping the Yeast N-Glycoproteome by Mass Spectrometry (MS). *Mol. Cell Proteomics*
32 **2014**, 13 (6), 1563–1572.
- 33 (54) Langó, T.; Róna, G.; Hunyadi-Gulyás, É.; Turiák, L.; Varga, J.; Dobson, L.; Várady, G.;
34 Drahos, L.; Vértessy, B. G.; Medzihradszky, K. F.; Szakács, G.; Tusnády, G. E.
35 Identification of Extracellular Segments by Mass Spectrometry Improves Topology
36 Prediction of Transmembrane Proteins. *Scientific Reports* **2017**, 7 (1), 1–9.
- 37 (55) Elia, G. Biotinylation Reagents for the Study of Cell Surface Proteins. *PROTEOMICS*
38 **2008**, 8 (19), 4012–4024.
- 39 (56) Wang, J.; Yang, P.; Tang, B.; Sun, X.; Zhang, R.; Guo, C.; Gong, G.; Liu, Y.; Li, R.;
40 Zhang, L.; Dai, Y.; Li, N. Expression and Characterization of Bioactive Recombinant
41 Human α -Lactalbumin in the Milk of Transgenic Cloned Cows. *Journal of Dairy Science*
42 **2008**, 91 (12), 4466–4476.
- 43 (57) Sato, C.; Kim, J.-H.; Abe, Y.; Saito, K.; Yokoyama, S.; Kohda, D. *Characterization of the*
44 *IV-Oligosaccharides Attached to the Atypical Asn-X-Cys Sequence of Recombinant*
45 *Human Epidermal Growth Factor Receptor*, 2000; Vol. 127.
- 46 (58) Zielinska, D. F.; Gnad, F.; Wiśniewski, J. R.; Mann, M. Precision Mapping of an in Vivo N-
47 Glycoproteome Reveals Rigid Topological and Sequence Constraints. *Cell* **2010**, 141 (5),
48 897–907.
- 49 (59) BOBBITT, J. M. Periodate Oxidation of Carbohydrates. *Advances in carbohydrate*
50 *chemistry* **1956**, 48 (11), 1–41.

- 1 (60) Muntel, J.; Boswell, S. A.; Tang, S.; Ahmed, S.; Wapinski, I.; Foley, G.; Steen, H.;
2 Springer, M. Abundance-Based Classifier for the Prediction of Mass Spectrometric
3 Peptide Detectability upon Enrichment (PPA). *Molecular & cellular proteomics*: MCP
4 **2015**, *14* (2), 430–440.
- 5 (61) Fusaro, V. A.; Mani, D. R.; Mesirov, J. P.; Carr, S. A. Prediction of High-Responding
6 Peptides for Targeted Protein Assays by Mass Spectrometry. *Nature biotechnology* **2009**,
7 *27* (2), 190–198.
- 8 (62) Eysers, C. E.; Lawless, C.; Wedge, D. C.; Lau, K. W.; Gaskell, S. J.; Hubbard, S. J.
9 CONSeQuence: Prediction of Reference Peptides for Absolute Quantitative Proteomics
10 Using Consensus Machine Learning Approaches. *Molecular and Cellular Proteomics*
11 **2011**, *10* (11).
- 12 (63) Searle, B. C.; Egertson, J. D.; Bollinger, J. G.; Stergachis, A. B.; MacCoss, M. J. Using
13 Data Independent Acquisition (DIA) to Model High-Responding Peptides for Targeted
14 Proteomics Experiments. *Molecular and Cellular Proteomics* **2015**, *14* (9), 2331–2340.
15

UniProt

Protein Sequence



MS suitability

in silico protein digest,
filter by peptide length
and ion m/z

Motif localization

- identify motif sites in
protein and peptides,
annotate by topology

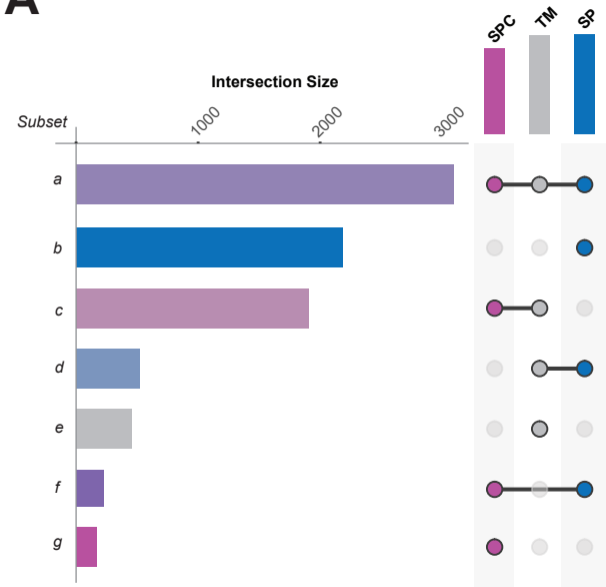
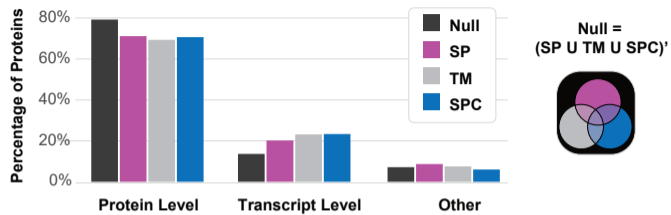
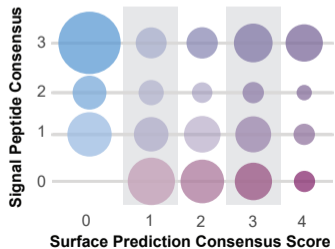
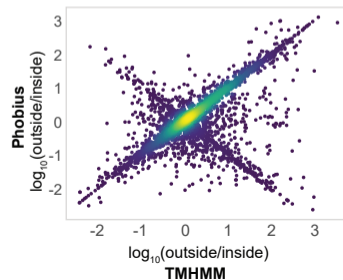
Protein annotation

- integrate individual
predictions to classify
proteins into subsets

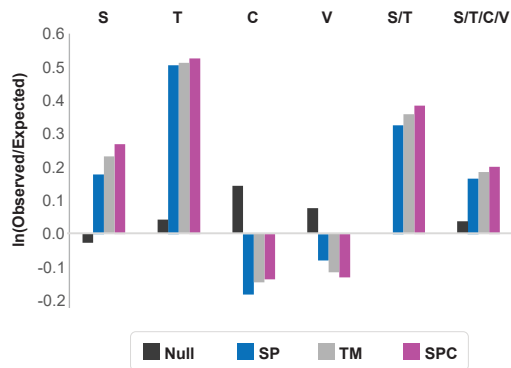


CIRFESS

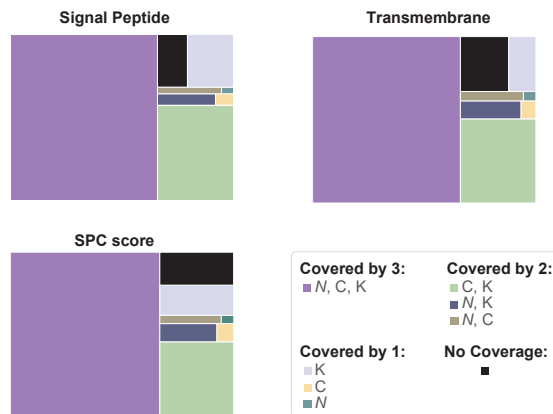
Compiled Interactive Resource For Extracellular and Surface Studies

A**B****C****D**

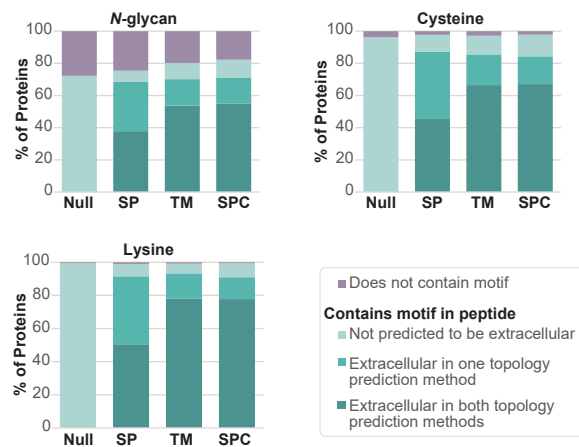
A



C



B



D

