

G OPEN ACCESS

Citation: Srivastava H, Lippincott MJ, Currie J, Canfield R, Lam MPY, Lau E (2022) Protein prediction models support widespread post-transcriptional regulation of protein abundance by interacting partners. PLoS Comput Biol 18(11): e1010702. https://doi.org/10.1371/journal.pcbi.1010702

Editor: Lukas Käll, KTH Royal Institute of Technology: Kungliga Tekniska Hogskolan, SWEDEN

Received: May 11, 2022

Accepted: November 1, 2022

Published: November 10, 2022

Copyright: © 2022 Srivastava et al. This is an open access article distributed under the terms of the Creative Commons Attribution License, which permits unrestricted use, distribution, and reproduction in any medium, provided the original author and source are credited.

Data Availability Statement: Code to train machine learning models and evaluate prediction performance have been uploaded to GitHub at https://github.com/Lau-Lab/CPTAC_Protein. Code for making figures can be found under notebook/ PLOS_MakeFig*.qmd. Supplementary Tables can be found on figshare at https://doi.org/10.6084/m9.figshare.19330541.

RESEARCH ARTICLE

Protein prediction models support widespread post-transcriptional regulation of protein abundance by interacting partners

Himangi Srivastava ^{1,2}, Michael J. Lippincott¹, Jordan Currie ¹, Robert Canfield ¹, Maggie P. Y. Lam^{1,2,3}, Edward Lau ^{1,2}*

- 1 Department of Medicine/Cardiology, University of Colorado School of Medicine, Aurora, Colorado, United States of America, 2 Consortium for Fibrosis Research and Translation, University of Colorado School of Medicine, Aurora, Colorado, United States of America, 3 Department of Biochemistry and Molecular Genetics, University of Colorado School of Medicine, Aurora, Colorado, United States of America
- * edward.lau@cuanschutz.edu

Abstract

Protein and mRNA levels correlate only moderately. The availability of proteogenomics data sets with protein and transcript measurements from matching samples is providing new opportunities to assess the degree to which protein levels in a system can be predicted from mRNA information. Here we examined the contributions of input features in protein abundance prediction models. Using large proteogenomics data from 8 cancer types within the Clinical Proteomic Tumor Analysis Consortium (CPTAC) data set, we trained models to predict the abundance of over 13,000 proteins using matching transcriptome data from up to 958 tumor or normal adjacent tissue samples each, and compared predictive performances across algorithms, data set sizes, and input features. Over one-third of proteins (4,648) showed relatively poor predictability (elastic net r < 0.3) from their cognate transcripts. Moreover, we found widespread occurrences where the abundance of a protein is considerably less well explained by its own cognate transcript level than that of one or more trans locus transcripts. The incorporation of additional trans-locus transcript abundance data as input features increasingly improved the ability to predict sample protein abundance. Transcripts that contribute to non-cognate protein abundance primarily involve those encoding known or predicted interaction partners of the protein of interest, including not only large multi-protein complexes as previously shown, but also small stable complexes in the proteome with only one or few stable interacting partners. Network analysis further shows a complex proteome-wide interdependency of protein abundance on the transcript levels of multiple interacting partners. The predictive model analysis here therefore supports that protein-protein interaction including in small protein complexes exert post-transcriptional influence on proteome compositions more broadly than previously recognized. Moreover, the results suggest mRNA and protein co-expression analysis may have utility for finding gene interactions and predicting expression changes in biological systems.

Funding: This work was supported in part by the U. S. National Institutes of Health (NIH; https://nih.gov) Office of the Director award R03-0D032666; NIH/NHLBI award R00-HL144829 to EL; and NIH/NHLBI awards R00-HL127302, R01-HL141278, and the Consortium for Fibrosis Research & Translation funds at the University of Colorado (http://cfret.org) to ML. The funders had no role in study design, data collection and analysis, decision to publish, or preparation of the manuscript.

Competing interests: The authors have declared that no competing interests exist.

Author summary

The abundance of mRNA is often measured as a surrogate variable of protein levels, but how well the mRNA level of different genes correlate with their protein across samples remains incompletely understood. Here we trained machine learning models over large RNA sequencing and mass spectrometry data from up to 8 cancer types in the CPTAC data sets to evaluate how well protein level variances across samples can be predicted from their transcripts. Despite voluminous data, up to one-third of genes shows poor mRNAprotein correlation suggesting their protein abundance is not primarily regulated from cognate transcripts. The inclusion of mRNA level information from protein interaction partners into the prediction models substantially improved prediction performance for a subset of genes, suggesting their protein abundance may be primarily regulated post-transcriptionally through protein-protein interactions. Notably, these proteins involve not only subunits of large multi-protein complexes such as the ribosome as previously suspected, but many proteins that form stable interactions with one or few other partners, including the propionyl-CoA carboxylase, mitochondrial calcium uniporter, calcineurin, and others. The results add to emerging evidence of independent regulation of protein levels from their cognate transcripts and suggest avenues to improve the interpretation of transcriptomics data.

Introduction

Mounting evidence now shows that protein levels correlate imperfectly with the levels of their cognate transcripts [1-3]. More specifically, although a robust trend exists over the log scale between protein and mRNA measurements across genes, genewise correlation between proteins and their transcripts is much poorer across observations (samples, tissues, cell types, or subjects). This has been taken to indicate that while abundant proteins have abundant transcripts, transcript variance within a group of samples does not necessarily predict or signify corresponding protein changes [4]. Multiple factors are known to contribute to this non-correlation. Technical variations are often cited as a substantial source of non-correlation, as transcriptomics and proteomics measurements carry different sources of error and proteins with lower baseline variance in mass spectrometry have been shown to be better predicted by their transcripts [5]. Nevertheless, a substantial portion of protein variance remains unexplained and is likely attributable to biological and biophysical regulations. It has been well recognized that large multi-protein complexes could invoke a buffer effect on protein levels [6,7], as a multimeric complex only fully folds and functions when all subunits are present, any induction of the transcript for a single subunit would not per se lead to additional complexes, and resulting supernumerary proteins are thought to be quickly degraded [6,8]. Lastly, numerous post-transcriptional and post-translational mechanisms are known to modulate protein levels such as the gene- and context-dependent translation rates of mRNAs [9,10], the differential half-life and temporal distributions between mRNAs and proteins [2,11], and proteolytic degrading of translated proteins in the cell.

The emergence of large-scale proteogenomics data from matching samples has created new opportunities to revisit protein-level predictions from transcriptomics data. The abundance of a protein may be the function of one or more transcripts. Most notably, available data sets from the Gene Tissue Expression (GTEx) project [7] and the Clinical Proteomic Tumor Analysis Consortium (CPTAC) [12], have spurred the use of machine learning approaches to evaluate how well one can predict protein level variance given a set of transcriptomics data,

with the goal of developing strategies and algorithms that can boost the performance of protein level predictions. This culminated in a community based effort in the CPTAC Proteogenomics Dream Challenge Task 2, which tasked participants with predicting protein abundances from mRNA and genetic data from CPTAC ovarian and breast cancer samples [12]. The results suggest that protein level prediction remains a challenging and not fully resolved problem, as many community-submitted models did not improve substantially the baseline model, which is an elastic net taking into account all mRNA features available and has a median Pearson's correlation coefficient (r) of 0.47 for ovarian cancer. Nevertheless, general lessons have emerged from the top performing models; for instance: (i) ensemble methods generally performed well [12-14]; (ii) combining observations from the ovarian and breast cancer datasets to borrow information from each other led to improved predictions [12]; and (iii) judicious feature pre-selection based on prior biological knowledge such as protein-protein interactions improved prediction performance [12,14]. Notwithstanding these general observations, the current literature reflects that much remains to be learned about the relationship of mRNA and protein regulations in different genes and whether there are fundamental limits to how well mRNA abundance reflects that of their protein counterpart. This problem has several practical importances. Proteins carry out the majority of biological processes and hence are arguably the most relevant molecules to biological states. Despite rapid advances in proteomics techniques, bulk and single-cell RNA sequencing remain the most commonly used methods to interrogate gene expression status on a large scale and will likely remain so in the foreseeable future. Transcriptomics experiments often operate on the implicit assumption that identified differential regulation exert their biological effects via their cognate proteins, hence it is important to better understand the relationships between protein and mRNA levels to aid in data interpretation and determining potential protein level changes given a set of transcriptomics data. Alternatively, knowing the genewise difference in how well a gene's transcript can predict its protein counterpart may be useful for filtering and prioritizing biologically relevant transcript signatures [15].

Here we revisit the predictability of protein levels from transcriptomics data. Since the time of the Dream Challenge, considerably more proteogenomics data have been made publicly available which increases the number of observations available for modeling training, as well as the number of proteins for which there is mass spectrometry information available. Individual CPTAC cancer studies have analyzed the protein and mRNA correlation in individual tumors and normal adjacent tissues and nominated specific pathways whose correlations are particularly poor. Additional re-analysis and meta-analysis studies have outlined the distribution of prediction performances across algorithms, and generally conclude there is some statistical enrichment of biological processes or protein features among proteins that are poorly predicted by their own transcript level, e.g., metabolic and essential proteins or proteins belonging to complexes [12,16]. Nevertheless, a granular analysis remains unrealized in the literature that interrogates the identity and regulatory modality of individual proteins in depth. Accordingly, our goals here are to (1) evaluate how the increasing data size from combining CPTAC tumor data sets affects the performance of prediction algorithms and feature selection strategies; and (2) interpret prediction models to assess the importance of transcript features in individual protein abundance regulation. The results suggest that the incorporation of transcript level information from protein interacting partners played a substantial role in predicting protein levels, and moreover, there are widespread instances in the proteome where the abundance of a protein correlates primarily with a trans locus transcript than its own cognate transcript, which has implications for gene expression profiling studies.

Methods

Data retrieval and processing

Gene expression data were obtained from public data from the CPTAC project and included data from 8 cancer types: ovarian cancer (OV) [17], breast cancer (BR) [18], endometrial carcinoma (EN) [19], colorectal cancer (CO) [20], lung adenocarcinoma (LUAD) [21], clear cell renal carcinoma (CCRCC) [22], glioblastoma (GB) [23], and lung squamous cell carcinoma (LSCC) [24]. The cumulative inclusions of each cancer type in the order above are sequentially referred to as CPTAC 2 to CPTAC 8 in the manuscript, such that CPTAC 2 refers to the union of ovarian and breast cancer (OV + BR); CPTAC_3 refers the union of ovarian, breast, and endometrial cancer (OV + BR + EN); and so on. The mRNA and protein level expression data from the CPTAC cancer types was retrieved using the cptac package v.0.9.7 [25] in Python 3.9. Each column of the quantitative measurement of the transcriptomics data acted as an independent variable or feature variable whereas the normalized quantitative measurement of a particular protein of interest acted as the single dependent or target variable in the protein model. Retrieved mRNA level gene expression data are standardized using the scikit-learn simple scaler. The proteomics data were likewise downloaded using the cptac package as presented in the data, and were stable isotope labeled relative quantitative mass spectrometry data presented as normalized log ratios across samples as in the original studies. All tumor samples were labeled using Thermo tandem mass tag (TMT) 10- or 11- plex isobaric tags for MS2 quantification, with the exception of the ovarian cancer data, which were labeled with Sciex iTRAQ isobaric tags for MS2 quantification, and the colon cancer data, which contained both label-free and TMT quantifications. The retrieved log ratios across samples were not further transformed.

For each protein for which predictions are to be made, we retrieved five separate feature sets:

- 1. **Single:** Using only the single transcript coding for the protein of interest for model training, then running the pipeline to train a model for each protein separately.
- 2. **CORUM:** Using the transcripts of all proteins belonging to the same protein complex as the protein of interest, if any, in CORUM v.3.0 [26], where protein complexes are defined as two or more proteins that interact physically in a quaternary structure. These transcripts then act as independent variables (features) to predict the target variable (protein of interest). The pipeline is then run to train a model for each protein separately.
- 3. **STRING 800:** Using the transcripts of interacting partners of the protein of interest as input features. Interacting partners are retrieved from STRING v.11 [27], which documents functional associations including physical interactions, genetic interactions, co-expression, co-occurrence, and other associations. The STRING combined score represents the overall likelihood of interactions. Interacting pairs with a STRING combined score of 800 or above (high-confidence) are included. The pipeline is then run to train a model for each protein separately.
- 4. **STRING 200:** As above, except that interacting pairs with a STRING combined score of 200 or above (low- to high-confidence) are included. The pipeline is then run to train a model for each protein separately.
- 5. **Transcriptome:** A transcriptome-wide model where all qualifying transcripts in the data set are included as features, prior to the removal of low-variance features. The pipeline is then run to train a model for each protein separately.

Model training and evaluation

For each feature set, missing values for feature variables are imputed using median imputation, followed by the removal of features with variance of 0.2 or below. Models are not trained for proteins with fewer than 50 empirical observations. The data are then split 80:20 into training and test sets. No imputation or additional standardization was performed on the proteomics data. The input and target data from the training set are then used to train a model using either linear regression, elastic net with 5-fold cross validation, random forest regressor, or gradient boosting regressor in scikit-learn v.1.0 [28], with the following specified parameters: random forest number of estimators: 500, criterion: squared error, max depth: 4; elastic net cross validation L1 ratio: 0.1, 0.5, 0.9, 0.95, cv: 5, tolerance: 1e-3, max iterations: 2000; gradient boosting: n estimators: 1000, max depth: 3, subsample: 0.5, min samples splot: 5, learning rate: 0.025. The trained models are saved as individual objects which include the predicted protein levels in the training set as well as the contribution of each feature to the overall prediction (coefficients in linear regression and elastic net; feature importance and trees in random forest and gradient boosting regressors), and are applied to predict protein levels in the test set data. Further interpretation of feature importance was performed using Shapley values with the aid of the shap package v.0.40.0 [29], or with the Boruta algorithm using the Boruta_Py package v.0.3 [30].

To evaluate model performance, the Pearson's correlation coefficients between predicted (\hat{y}) and actual (y) protein values are calculated individually for the train and test set data for each protein model using the numpy corrcoef function, or defaulted to 0 if the standard deviation of predicted values is 0. Goodness-of-fit (R^2) is calculated using scikit-learn.metrics. r2_score function. Normalized root mean square errors (NRMSE) are calculated as follows:

$$NRMSE = \sqrt{\frac{1}{n} \sum_{i=0}^{n-1} (y_i - \hat{y}_i)^2 * \frac{1}{(max(y) - min(y))}}$$

Network construction and analysis

To construct the protein regulation network, we first constructed an overall graph from a subset of transcript-protein relationships using the networkx package [31] in Python 3.9. Qualifying proteins are those whose prediction increased with the inclusion of more features, such that in the STRING or CORUM feature set, the correlation coefficient between the elastic net predicted and actual protein levels is greater than the single feature elastic nets by at least 0.25. Transcripts that contribute to the prediction of these proteins are therefore included if their elastic net coefficient a and random forest feature importance b are greater than a certain threshold, which was set to a = 0.05, b = 0.05 for the CORUM feature set and a = 0.2, b = 0.05 for the STRING feature set.

In the data, the list of proteins and their corresponding transcriptomics features act as the nodes and the interaction pattern is a function of the correlation value and feature importance between a protein and the transcriptomics. The direction of the edge therefore flows from transcriptomics to protein. All the text files for each protein obtained from the computational pipeline are compiled into a n overall data frame. This data frame is then converted into the form of a network data frame. In this data frame, all the protein becomes the target node, and the transcriptomics corresponding to the protein becomes the source node. Additional column of this data frame consists of the weight between a protein and its predicting transcripts, which was calculated using the sum of the [–1,1] clipped elastic net coefficient and the descending percentile rank of the random forest feature importance. An overall directed graph G(V, E) is then constructed from the edge lists with proteins as targets and their predicting transcripts as

sources, such that it is a function of vertices V which depict protein and/or transcript nodes and edges E which flow in the direction from the transcript sources to the protein targets. Using the built-in functions in networkx, this overall directed graph is divided into weakly connected components, which resulted in individual subgraphs such that for the overall network G = (V, E) where V is the vertices and E is the edge then a subgraph of G = (V, E) is a graph S = (V', E') where vertex set $V' \subseteq V$ and edge set $E' \subseteq E$ connects only nodes of V'. We then used Cytoscape v.3.9 [32] to visualize and perform additional topology analysis and functional annotations of each subgraph. HiDef persistent community detection [33] was performed with the aid of the CyCommunityDetection Cytoscape plugin [34] using the Leiden algorithm [35]. Functional enrichment of individual communities was performed using CyCommunityDetection with the g:Profiler web server [36]. Hub nodes were identified with the aid of the CytoNCA [37] Cytoscape plugin using the betweenness centrality algorithm with the top 10% of proteins with the highest centrality defined as hubs. Functional enrichment analysis on the hub proteins was performed with the aid of the stringApp [38] Cytoscape plugin using default settings and 5% false discovery rate cutoff.

Additional data analysis

Additional data analysis, statistics, and visualization were performed in R v.4.1.1 with the aid of the ReactomePA [39], clusterProfiler [40], and circlize [41] packages; and in Python 3.9 with the aid of the seaborn [42] package.

Results

Feature selection improves the prediction of protein abundance from transcriptome data

We first evaluated whether increasing proteogenomics data depth would impact optimal feature selection strategies and algorithms in prediction protein levels. To do so, we retrieved data from up to 8 tumor types from the CPTAC data set and processed the data to collage a data table with matching protein and transcript data from each tumor sample. The data sets were ordered such that breast and ovarian cancers as in the DREAM challenge were combined first, then other cancer types were added according to the order of their availability (see Methods). We trained models to predict the target variable (the normalized labeled mass spectrometry measured abundance of a particular protein) from various input features (transcriptomics data from the matching samples) using either a multiple linear regression, elastic net, and random forest regressor in scikit-learn. To compare the performance of prior knowledge based feature selection, we further compared supplying the models with input features based on (1) only RNA level data the corresponding transcript of the target variable protein (single feature) (2) self transcript plus the transcript of any protein within the same protein complex of the target variable protein in CORUM (CORUM feature); (3) self transcript plus the transcript of any proteins that are high-confidence protein interaction partners including physical interactions and other inferred association from STRINGdb (STRING 800 feature); (4) self transcript plus the transcript of any proteins that are low to high- confidence protein interaction partners from STRINGdb; (5) using all qualifying transcripts (transcriptome features).

We next compared the genewise prediction performance of the transcriptome using the test set correlation coefficient between predicted and actual mass spectrometry protein values, R², and NRMSE as described (Fig 1A and S1 Table). The elastic nets and random forest models performed better with more features and more data. In contrast, a multiple linear regression model performed little better with more included features and in fact failed to predict

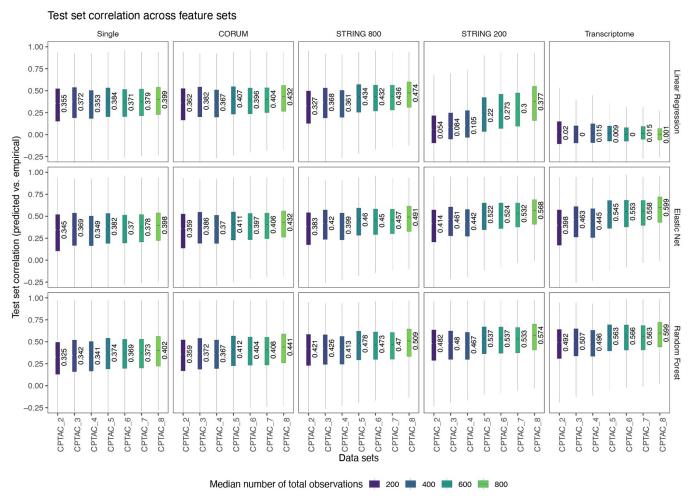


Fig 1. Genewise dispersion of protein predictability from transcriptome data. Box plots of test set correlation coefficients between the transcript-predicted and actual protein level for each protein are shown across five feature sets (column: single/self transcript, CORUM interactors, STRING 800 high-confidence associated proteins; STRING 200 low-confidence associated proteins, and all transcripts) and three algorithms (multiple linear regression, elastic net, and random forest). In each plot, the x axis denotes the number of additive CPTAC data sets used to train the models as described in Methods; box: interquartile range; whiskers: +/- 1.5 IQR; notch: SEM.

https://doi.org/10.1371/journal.pcbi.1010702.g001

protein levels when a large number of features are given. For the elastic net and random forest models, performance gains began to saturate with additional data sets (r: 0.388 and 0.402, respectively) but further increased when larger feature sets were given (r: 0.599 for both) (**Figs 1A and S1**). As a comparison, we also evaluated the performance of the models in single cancer data (i.e., in each data set alone) (**S2A Fig**) as well as combined the data sets in the order of descending performance (**S2B Fig**). The single cancer models suggest that a prediction performance of up to median r of 0.691 is achievable using the Transcriptome wide feature set, and 0.69 in the STRING 200 feature set (random forest, CCRCC), albeit in a single cancer type only. In both the single cancer type and reordered combination comparisons, the observation remains true in all cases that the inclusion of protein interactor features further increased the prediction performance of both the elastic net and random forest algorithms. Hence, substantial prediction improvements resulted from the inclusion of other transcripts (number of features) over the gain in data set sizes (number of observations) or the use of algorithms (random forests) that can account for non-linearity. This is corroborated when considering the feature sets across the largest data set collection used (CPTAC_8 from 8 cancer types).

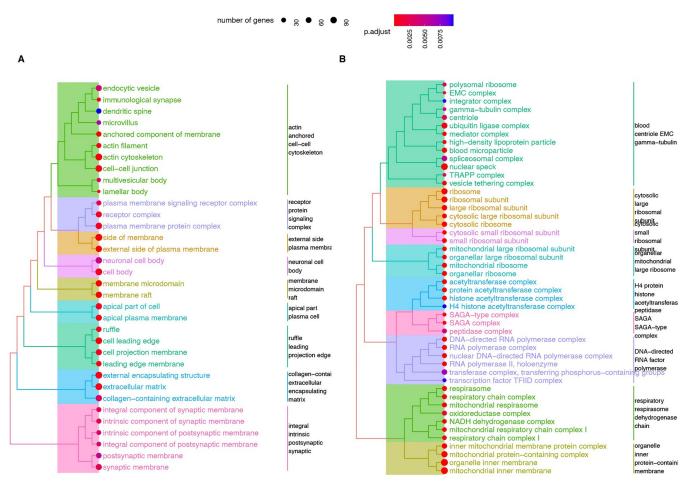


Fig 2. Pathway enrichment of proteins with good and poor predictability. A. Tree plots showing the clustering and relationships of gene ontology terms that are significantly enriched among proteins whose abundances are well predicted by their own transcripts ($r \ge 0.6$). B. Tree plots of terms enriched among proteins whose abundances are poorly predicted by their own transcripts ($r \le 0.3$).

https://doi.org/10.1371/journal.pcbi.1010702.g002

Model performance continued to increase as more features were added in (median r from 0.398 to 0.599), and the number of non-predictors and negative predictors decreased while overall dispersion of predictability also decreased.

From the baseline self-transcript model in the CPTAC_8 data set containing 8 cancer types, we found that strong protein predictors ($r \ge 0.6$; 2,008 genes) are enriched in genes participating in membrane proteins and cell junctions (**Fig 2A and S2 Table**), whereas poor protein predictors ($r \le 0.3$; 4,648 genes) are enriched in multiple large multi-protein complexes (**Fig 2B and S3 Table**). Hence the analysis of CPTAC_8 data set confirms prior observations that protein complex membership presents a major source for transcript-protein non-correlation.

We next compared the genewise performances between the STRING feature elastic net model and the self transcript model to investigate the observed performance gain at a more granular level. Considering the 13,239 genes with protein-level prediction and with at least one STRING interactor, on average, each gene in the STRING model has a median of 273 [158–457] features, compared to 1 feature in the self transcript model. Incorporation of STRING features led to an average increase in test set correlation coefficients of 0.16 [IQR: 0.07–0.27] (S4 Table). There is a significant positive correlation between the number of features and prediction improvements (Pearson's r 0.21, P: 8.2e–131) (S3A Fig), a relationship which persisted

when only genes with 10 or more interactors were included S3B Fig). A similar observation is seen in the CORUM feature set where the number of transcript features (i.e., CORUM complex interactors) of a particular protein is positively correlated with improvement in predictability of protein abundance (S3C and S3D Fig). A functional enrichment analysis of proteins with strong improvements in prediction (difference in STRING vs. self-feature in test set predicted-actual correlation coefficient (Δr) \geq 0.25) showed a strong enrichment in large multiprotein complexes including the ribosome, mitochondrial ribosome, RNA polymerase, and spliceosome (S5 Table) similar to the results from poor predictors. This corroborates that the non-correlation between transcripts and proteins of multi-protein complexes levels can be partially rescued by considering the information of interacting partner transcripts, which could participate in post-transcriptional regulations such as by being the stoichiometrically limiting transcript or crucial assembly parts, such that the transcript level of one subunit can regulate the protein level of other subunits.

Overall, these results are consistent with information about protein level residing both in the transcripts of the genes encoding the protein as well as the protein association partners. Moreover, taking into consideration the number of features in the STRING 800 and STRING 200 sets (median feature sizes of 10 [IQR: 2–48] and 266 [IQR: 149–451], respectively) compared to the size of the whole transcriptome feature set, the results indicate that a substantial portion of performance gains over the single-transcript feature set is already achievable from relatively few selected prior features.

Proteins whose abundances are predicted by non-cognate transcripts are common

We next examined more closely the underlying causes behind the performance gains of protein predictions in the CORUM and STRING feature sets. To do so, we examined the genes whose protein prediction performance from transcripts increased substantially ($\Delta r > 0.25$) after the incorporation of additional transcript features. In total, we observed 484 such proteins in the CORUM feature sets, and 3,272 proteins from the STRING data set, representing over 24% of all examined proteins. These numbers increase further when a less conservative Δr threshold of ≥ 0.15 is used (6,123 proteins in the STRING dataset, 946 in CORUM), altogether suggesting non-cognate transcript contributions to protein level are common at the proteome level. Upon inspecting the model feature coefficients (elastic nets) and feature importance (random forests), we observed that a considerable portion of these proteins with improved prediction are associated with (1) poor contribution from the self-transcript, and (2) a substantial contribution from primarily a few non-cognate (i.e., trans locus) transcripts coding for other proteins (see below). In other words, although model performance continued to increase with increased feature set sizes, the contributions of trans locus transcript features to overall prediction performance is unevenly distributed, and are therefore attributable to a few high contribution genes rather than a simple scaling with feature size. Notably, despite the STRING feature set being substantially larger than the CORUM feature set, substantial contributions from trans locus transcripts in the STRING feature set primarily involve transcripts encoding proteins that form part of the same CORUM complex as the protein of interest itself, suggesting that stable complex memberships play an outsized role in determining protein levels.

To illustrate the poor predictive power of cognate transcripts on some proteins, we next interrogated a subset of proteins where a non-self transcript has an outsized effect on the protein level of the proteins in the CORUM and STRING feature sets: Propionyl-CoA Carboxylase Subunit Beta (PCCB), C-X9-C Motif Containing 1 (CMC1), Proteasome Assembly

Chaperone 2 (PSMG2), SMCR8-C9orf72 Complex Subunit (SMCR8), Mitochondrial Calcium Uptake 2 (MICU2), and Protein Phosphatase 3 Regulatory Subunit B, Alpha (PPP3R1). PCCB forms the propionyl-CoA carboxylase complex with two protein members, which breaks down certain amino acids in the cell (Fig 3A). PCCB showed poor prediction by their own self-transcripts (r:-0.008) but experienced high gains in test set correlation upon incorporating additional transcript features in the CORUM (r: 0.483) and STRING (r: 0.692) models. CMC1 is an assembly factor that forms an early intermediate of the cytochrome c oxidase complex in the mitochondrion with 14 documented protein members in CORUM (Fig 3B). Likewise, CMC1 showed only moderate prediction by their own self-transcripts (r: 0.274, respectively) but experienced high gains in test set correlation upon incorporating additional transcript features in the CORUM (r: 0.615) and STRING (r: 0.564) models.

We then considered the transcriptomics and proteomics data distributions of the features that had highest coefficients in the elastic net model from the CORUM feature set for PCCB and CMC1. For PCCB, most protein-level predictions are recovered when the feature set includes PCCA, which together with PCCB forms the stable propionyl-CoA carboxylase enzyme that consists of six copies of PCCB and six copies of PCCA each. There is a remarkably strong correlation with its propionyl-CoA carboxylase complex interacting partner PCCA at the protein level (r: 0.957), but this correlation is almost entirely absent at the transcript level (r: -0.059). This is partially explainable by the observation that PCCA transcript has a higher level of variability and is also strongly correlated with the PCCB protein level. The total protein level of PCCB is therefore primarily driven by PCCA rather than PCCB transcripts (Fig 4A).

In the case of CMC1, we found that MT-CO1 has the highest contribution to CMC1 protein level (Fig 4B). MT-CO1 is the mitochondrial genome encoded subunit of complex IV that is thought to act as the nascent scaffold around which the complex is assembled. CMC1 binds to MT-CO1 during the formation of the early module known as MITRAC and is subsequently released during assembly [43,44]. CMC1 transcript and protein levels are moderately correlated with one another (r: 0.228). Again, CMC1 and MT-CO1 show poor co-expression at the transcript (r: -0.017) level and an improved correlation at the protein (r: 0.339) level, whereas MT-CO1 transcript is robustly correlated with the CMC1 protein level (r: 0.583), which suggests the possibility that proteogenomic co-expression may reveal additional functionally related protein pairs. Closer inspection suggests the possibility of a non-linear relationship between MT-CO1 transcript and CMC1 protein where upon reaching a plateau, further increases in MT-CO1 transcript does not correlate to further increase in CMC1 protein levels, suggesting other transcripts may also contribute to CMC1 protein levels. Consistent with this, other complex IV subunits also have non-negligible coefficients (in the elastic nets) or feature importance (in the random forests) in the CMC1 model. When the feature set expanded to include all STRING proteins, MT-CO1 remained a high contributor in the random forest but not the elastic net model, suggesting the elastic nets may be more susceptible to collinear features than the random forests.

Two other examples of small stable complexes with interdependent protein and mRNA correlation are shown in Fig 5. Proteasome Assembly Chaperone 2 (PSMG2) forms the heterodimer PAC1-PAC2 complex (2 protein members in CORUM) along with Proteasome Assembly Chaperone 1 (PSMG1), which binds with 20S proteasome precursors and acts as a scaffold to promote proteasome ring assembly while preventing aberrant dimerization of 20S proteasome α rings [45]. PSMG2 and PSMG1 are co-expressed strongly at the protein level (r: 0.713) but only moderately at the transcript level (0.222) (Fig 5A). PSMG2 protein has a higher correlation with PSMG1 transcript (r: 0.472) than its own (r: 0.158). From GTEx v8 data, PSMG2 is expressed more highly than PSMG1 across tissues (TPM ~40 vs. 15) [46], which suggests PSMG1 may act as a limiting factor in heterodimer formation. PSMG1 knockout in mice has

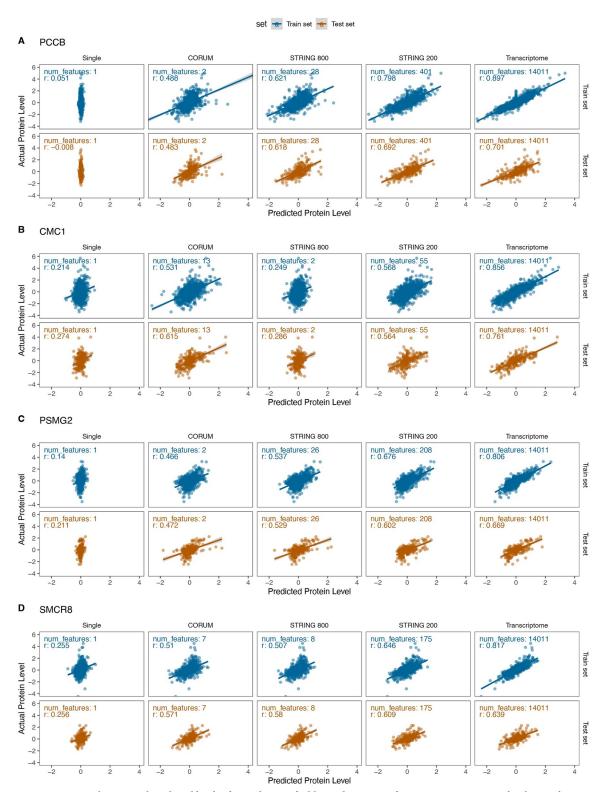


Fig 3. Proteins with improved predicted levels after inclusion of additional transcript features. Four proteins with substantial predictability from transcriptome data upon the inclusion of additional features are shown: **A.** PCCB, **B.** CMC1, **C.** PSMG2, **D.** SMCR8. For each protein, the transcript-trained prediction of protein level is plotted on the x axis and the actual protein level is plotted on the y axis. The lack of variance in predicted protein levels from the self-transcript model is due to the regularization of the elastic net model, and corresponds to a lack of correlation between PCCB mRNA and protein (see Fig 4). Blue: train set, brown: test set. Columns denote

the transcript feature set used to train the model. The number of features used to train the model in each feature set is shown inside each plot. r: Correlation coefficient.

https://doi.org/10.1371/journal.pcbi.1010702.g003

been shown to decrease PSMG2 at the protein level, consistent with our finding that PSMG1 transcript levels influence PSMG2 protein abundance [47].

In another example, Smith-Magenis Syndrome Chromosomal Region Candidate Gene 8 (SMCR8) forms the heterotrimer C9orf72-SMCR8-WDR41 complex (3 protein members in CORUM) together with the C9orf72-SMCR8 Complex Subunit (C9orf72) and WD Repeat Domain 41 (WDR41) [48]. C9orf72-SMCR8 is known to dimerize prior to binding with WDR41 to form the heterotrimer. The C9orf72-SMCR8 complex modulates autophagy, likely

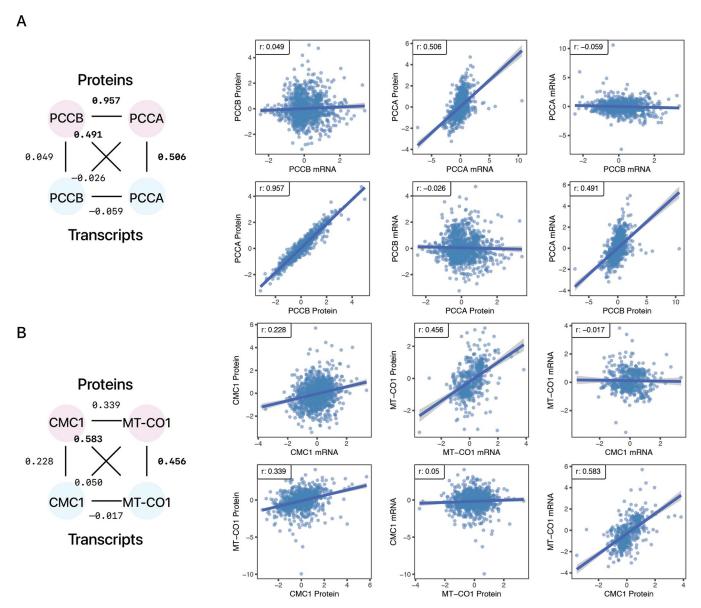


Fig 4. mRNA-Protein correlations of PCCB and CMC1 with functionally associated proteins. Two examples of proteins whose abundance is better explained by another transcript are shown. A. PCCB protein level is predicted by PCCA transcript but not its own transcript. B. CMC1 protein level is explained by MT-CO1 transcript level but not its own transcript. Substantial correlations across transcripts and proteins (\geq 0.4) are bolded.

https://doi.org/10.1371/journal.pcbi.1010702.g004

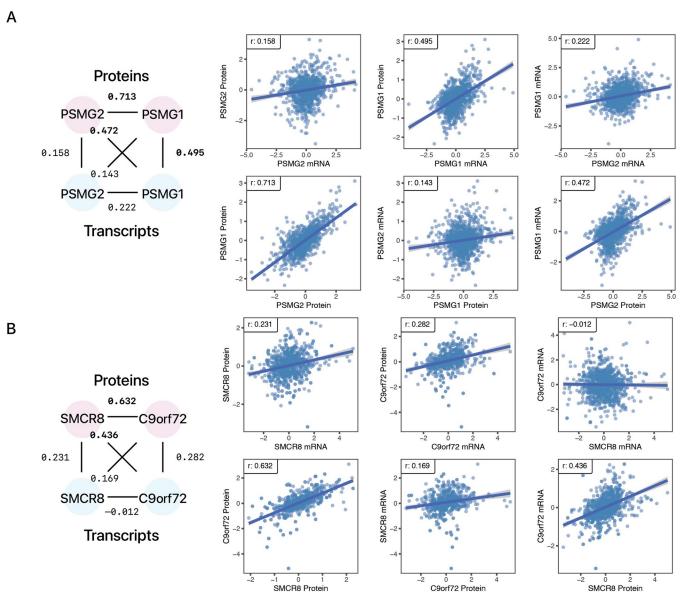


Fig 5. mRNA-Protein correlations of PSMG2 and SMCR8 with functionally associated proteins. Two examples of proteins whose abundance is better explained by another transcript are shown. A. PSMG2 protein level is predicted by PSMG1 transcript but not its own transcript. B. SMCR8 protein level is explained by C9orf72 transcript level but not its own transcript. Substantial correlations across transcripts and proteins (\geq 0.4) are bolded.

https://doi.org/10.1371/journal.pcbi.1010702.g005

by modulating the maturation of autophagosomes [49]. Mutations in the genes encoding the C9orf72-SMCR8 Complex are broadly implicated in diseases including amyotrophic lateral sclerosis and frontotemporal dementia. As in the other cases, we observed a co-expression of SMCR8 and C9orf72 at the protein level (0.632) but not the transcript level (-0.012), and SMCR8 protein levels correlate more strongly with the C9orf72 transcript (r: 0.436) than the SMCR8 transcript (r:0.231) (Fig 5B). A previous study suggested C9orf72 may stabilize excess SMCR8 [50], providing one possible explanation for the mechanism behind the observations here.

S5 Fig shows two more examples concerning small complexes. The mitochondrial calcium uniporter complex (MCU complex) contains five protein members in CORUM (S5A Fig).

Regulatory subunits of the complex, including calcium uptake protein 2, mitochondrial (MICU2) sense calcium levels, regulating the uptake of calcium performed by the mitochondrial calcium uniporter (MCU) [51]. Heterologous overexpression of FLAG-tagged MCU has been shown to increase expression of MICU2 in culture [52]. This is interesting, given the observation here that at an endogenous level of MCU expression, MICU2 protein is more strongly correlated to the MCU transcript (r: 0.488) than its own (r: 0.058) (S5A Fig).

The final example highlighted here comes from the heterodimer calcineurin, which is made up of a catalytic subunit (calcineurin subunit A), with 3 isoforms encoded by 3 separate genes (PPP3CA, PPP3CB and PPP3CC), and a Ca²⁺ binding regulatory subunit (calcineurin subunit B), with 2 isoforms also encoded by 2 separate genes (PPP3R1 and PPP3R2) [53]. Calcineurin is a serine/threonine phosphatase which modulates many calcium dependent signaling pathways, and is the target of immunosuppressant drugs cyclosporin A and FK506 [54]. In the analysis here, calcineurin subunit B type 1 (PPP3R1) protein correlated poorly with PPP3R1 transcript levels (r: 0.205), but strongly correlates (r: 0.819) with the protein phosphatase 3 catalytic subunit alpha (PPP3CA) transcript—an isoform of calcineurin subunit A (S5B Fig). Together, these cases illustrate a general observation—in cases where the protein abundance of a gene is more highly correlated to the transcript level of a binding partner than to its cognate transcript, there appears to be stronger co-expression between the gene and its abundance driver at the protein level but not at the transcript level.

To corroborate these observations, we used the Boruta algorithm against the random forest CORUM models and confirmed each of the principal non-cognate transcript features is retained. We further considered Shapley values as an explanation of random forest and gradient boosting models, and likewise found that the interacting protein transcripts formed the top contributors in explaining the abundance of the protein of interest in each of the six cases above (S6 Fig). In some examined cases including PSMG2, the transcript coding for the protein of action is more abundant than that of its regulating interacting partner, which is compatible with a scenario where simple stoichiometric constraints control the protein level of the supernumerary subunit. It is however less clear whether this scenario applies to other examined cases including CMC1 and MT-CO1, where the MT-CO1 transcript is much more abundant than CMC1 owing to the multiple copies of mitochondrial genome in the tissue. Hence, how MT-CO1 regulates CMC1 protein level is not explainable directly from transcript stoichiometry alone and awaits further mechanistic studies. Finally, we also examined feature importance in a single data set of a recent CPTAC study that showed the highest performance in the single feature model (LSCC) (S7 Fig). Among the six highlighted proteins (PCCB, CMC1, PSMG2, SMCR8, MICU2, PPP3R1), five were best predicted by a non-self transcript, and five of the top trans locus predictors were conserved from the combined data set analysis with the exception of MT-CO1, which was not among the model features due to the number of shared observations required. Hence the observed trans locus predictors are conserved in a single data set model and are unlikely to be due to heterogeneity across data sets.

Table 1 summarizes the top 50 genes whose protein abundance becomes substantially more predictable upon the inclusion of CORUM features, along with their top protein level contributor transcripts. Notably, a number of these genes are found in multiple gene sets within the MSigDB C2 CGP gene set collection that are commonly used for functional enrichment analysis [55], whereas others are found to be significantly associated with diseases in the literature from bibliometric analysis [56,57], hence these genes are common and important to multiple areas of biomedical inquiries. We suggest the possibility that the transcripts of these genes may not present accurate proxy variables for their proteins should be considered in linking transcript level changes to downstream cellular physiology.

Table 1. Top 50 proteins whose abundance is under substantial influence from non-cognate transcripts. Columns 1 and 2: gene names. Columns 3 show the representative CORUM complex the gene of interest belongs to. Columns 4 and 5 denote the increase in prediction performance between elastic net single feature (self transcript) vs. CORUM feature sets. Column 6 shows the number of transcripts used to predict the protein level of the gene of interest in the CORUM feature set. Column 7 shows the top trans-locus contributor to the protein level of the gene of interest, ranked by absolute coefficients in the elastic net model. Proteins whose own transcripts are the top predictors are marked with (self). Column 8 denotes the number of MSigDB C2 CGP (chemical and genetic perturbation) gene sets in which the gene appears. Column 9 denotes the top significantly associated Disease Ontology term with the gene of interest in the literature.

Gene	Gene name	Representative CORUM complex name	Test set r, single feature	Test set r, CORUM features	# CORUM feature	Top non-self contributor	# MSigDB CGP sets	Most significant Disease Ontology association $(P < 0.05)$
APOL1	apolipoprotein L1	APOL1 complex B (APOL1, APOA1, HPR, FN1, IGHM)	0.118	0.435	5	APOA1	20	Kidney disease (P = 0.00011)
BORCS8	BLOC-1 related complex subunit 8	BORC complex	0.189	0.456	8	BORCS7	5	
CACNA1A	calcium voltage-gated channel subunit alpha1 A	G protein complex (CACNA1A, GNB1, GNG2)	0	0.305	3	GNG2	31	Familial hemiplegic migrane (P = 0.00032)
CAPZB	capping actin protein of muscle Z-line subunit beta	CAPZA3-CAPZB complex	0.032	0.719	22	WASHC2A	19	
CMC1	C-X9-C motif containing 1	Cytochrome c oxidase, mitochondrial	0.274	0.615	13	MT-CO1	10	Osteoarthritis (P = 0.002)
COL18A1	collagen type XVIII alpha 1 chain	ITGA5-ITGB1-CAL4A3 complex	0.056	0.359	3	ITGA5	82	
CPLX3	complexin 3	SNARE complex (VAMP2, SNAP25, STX1a, STX3, CPLX1, CPLX3, CPLX4)	0.478	0.515	6	SNAP25	6	Asphyxia neonatorum (P = 0.013)
DAG1	dystroglycan 1	UTM-SGCE-DAG1-CAV1-NOS3 complex	0.251	0.604	6	CAV3	33	Muscle tissue disease (P = 0.0023)
DLG4	discs large MAGUK scaffold protein 4	DLG4-DLGAP1-SHANK3 complex	0	0.51	6	FYN	18	Schizophrenia 4 (P = 0.036)
EPN1	epsin 1	RalBP1-CCNB1-AP2A-NUMB-EPN1 complex	0.098	0.453	5	CCNB1	6	
ESPL1	extra spindle pole bodies like 1, separase	ESPL1-CDC2 complex	-0.372	0.1	3	PTTG1	55	Spindle cell cancer (P = 0.03)
F7	coagulation factor VII	Factor-Xa-TFPI-factor-VIIa-tissue factor complex	0	0.479	4	F10	27	Factor VII deficiency (P = 0.0014)
GNB1	G protein subunit beta 1	G protein complex (BTK, GNG1, GNG2)	0.195	0.665	28	PTH1R	53	
GP1BA	glycoprotein Ib platelet subunit alpha	ITGA2b-ITGB3-CD9-GP1b-CD47 complex	0.262	0.523	6	VWF	12	Blood platelet disease (P = 0.0006)
GTF2F1	general transcription factor IIF subunit 1	RNA polymerase II complex, (CDK8 complex)	0.18	0.644	44	GTF2F2	17	
HSPA1A	heat shock protein family A (Hsp70) member 1A	HSP70-BAG5-PARK2 complex	0.182	0.48	42	HSPA8	96	Ischemia (P = 0.017)
HUS1B	HUS1 checkpoint clamp component B	9b-1b-1 complex	0	0.57	2	RAD1	0	Testicular Leydig cell tumor (P = 0.032)
IDH3B	isocitrate dehydrogenase (NAD (+)) 3 non-catalytic subunit beta	Isocitrate dehydrogenase [NAD], mitochondrial	0.158	0.431	3	IDH3A	37	
IRAK2	interleukin 1 receptor associated kinase 2	IRAK1-IRAK2 complex	0.227	0.484	2	(Self)	33	Lymphocytic colitis (P = 0.018)

(Continued)

Table 1. (Continued)

Gene	Gene name	Representative CORUM complex name	Test set r, single feature	Test set r, CORUM features	# CORUM feature	Top non-self contributor	# MSigDB CGP sets	Most significant Disease Ontology association (P < 0.05)
ITGA2B	integrin subunit alpha 2b	ITGA2b-ITGB3-CD9-GP1b-CD47 complex	0.204	0.48	12	TGM2	30	Blood platelet disease (P = 0.00038)
KCNQ2	potassium voltage- gated channel subfamily Q member 2	KCNQ2-KCNQ3 complex	0.007	0.331	2	KCNQ3	13	Epilepsy (P = 0.0012)
LAMA5	laminin subunit	ITGA6-ITGB4-LAMA5 complex	0.192	0.466	6	ITGB4	47	Placental infarction (P = 0.024)
LRP5	LDL receptor related protein 5	Norrin receptor complex	0.121	0.442	3	FZD4	26	
MPP3	MAGUK p55 scaffold protein 3	CADM1-4.1B-MPP3 complex	-0.069	0.685	14	ITGB2	20	
MRPL53	mitochondrial ribosomal protein L53	39S ribosomal subunit, mitochondrial	0.21	0.654	78	MRPL12	4	
NDUFAF1	NADH:ubiquinone oxidoreductase complex assembly factor 1	Ecsit complex (ECSIT, MT-CO2, GAPDH, TRAF6, NDUFAF1)	0.149	0.443	12	ECSIT	14	Protein deficiency (P = 0.0038)
NEDD8	NEDD8 ubiquitin like modifier	Ubiquitin E3 ligase (CDC34, NEDD8, BTRC, CUL1, SKP1A, RBX1)	0.087	0.343	7	SKP1	22	
NGDN	neuroguidin	AATF-NGDN-NOL10 complex	0.249	0.538	3	NOL10	9	
NPHP4	nephrocystin 4	DVL2-INVS-NPHP4-RPGRIP1L complex	0	0.36	8	NPHP1	8	Kidney disease $(P = 0.0011)$
PATJ	PATJ crumbs cell polarity complex component	CRB1-MPP5-INADL complex	0.27	0.562	11	(Self)	26	
РССВ	propionyl-CoA carboxylase subunit beta	Propionyl-CoA carboxylase	-0.008	0.483	2	PCCA	30	Propionic acidemia (P = 0.0048)
PHF21B	PHD finger protein 21B	LSD1 complex	0	0.342	13	HMG20A	7	Colon squamous cell carcinoma (P = 0.03)
PHKG1	phosphorylase kinase catalytic subunit gamma 1	Phosphorylase kinase complex	-0.629	-0.17	3	(Self)	8	Alzheimer's disease (P = 0.041)
PPP3R1	protein phosphatase 3 regulatory subunit B, alpha	Calcineurin-FKBP12 complex	0.263	0.544	3	PPP3CA	21	
PSMB5	proteasome 20S subunit beta 5	26S proteasome	0.022	0.53	37	PSME2	43	Multiple myeloma (P = 0.037)
PSMB6	proteasome 20S subunit beta 6	26S proteasome	0	0.53	37	PSME2	31	
PSMC3IP	PSMC3 interacting protein	TBPIP/HOP2-MND1 complex	0.294	0.551	2	MND1	42	Recurrent ovarian germ cell neoplasm (P = 0.017)
PSMG2	proteasome assembly chaperone 2	PAC1-PAC2 complex	0.211	0.472	2	PSMG1	13	Herpes simplex (P = 0.03)
REV1	REV1 DNA directed polymerase	Rev1-Rev3-Rev7-Polkappa complex	0	0.5	4	REV3L	14	Brucellosis (P = 0.0008)

(Continued)

Table 1. (Continued)

Gene	Gene name	Representative CORUM complex name	Test set r, single feature	Test set r, CORUM features	# CORUM feature	Top non-self contributor	# MSigDB CGP sets	Most significant Disease Ontology association $(P < 0.05)$
RPA3	replication protein A3	RPA complex	0.211	0.537	20	RPA2	45	Combined thymoma (P = 0.048)
RPGRIP1L	RPGRIP1 like	DVL2-INVS-NPHP4-RPGRIP1L complex	-0.051	0.289	4	INVS	6	N syndrome (P = 0.0079)
RPS29	ribosomal protein S29	Nop56p-associated pre-rRNA complex	0.06	0.455	120	TUBB1	17	
SERPINA1	serpin family A member 1	SERPINA1-CTSG complex SERPINA1-ELA2 complex	0.201	0.467	3	CTSG	73	Lung disease (P = 0.000008)
SERPIND1	serpin family D member 1	MLL-HCF complex	0	0.581	7	MEN1	17	Cervical cancer (P = 0.00094)
SMCR8	SMCR8-C9orf72 complex subunit	WDR41-(C9orf72-SMCR8)- (FIP200-ULK1-ATG13-ATG101) complex	0.256	0.571	7	C9orf72	9	Amyotrophic lateral sclerosis (P = 0.026)
TEN1	TEN1 subunit of CST complex	CST complex	0	0.383	3	STN1	5	Coats disease (P = 0.016)
TNIP2	TNFAIP3 interacting protein 2	TNF-alpha/NF-kappa B signaling complex	0.161	0.629	25	MAP3K8	17	Malignant histiocytic disease (P = 0.022)
TTR	transthyretin	TTR-RBP complex	0.084	0.414	2	RBP4	21	Amyloidosis (P = 0.0000081)
UBA52	ubiquitin A-52 residue ribosomal protein fusion product 1	60S ribosomal subunit, cytoplasmic Ribosome, cytoplasmic	0	0.537	79	RPS9	20	Iridocyclitis (P = 0.026)
YBX3	Y-box binding protein 3	H2AX complex II	0.141	0.424	5	NPM1	55	Lyme disease (P = 0.0012)

https://doi.org/10.1371/journal.pcbi.1010702.t001

Taken together, these results exemplify widespread interdependency of protein levels on trans locus transcripts, involving not only large megadalton protein complexes but many small complexes involved in diverse biological processes in the cell. The CORUM feature set considered here alone contains over 2,700 human complexes derived from CORUM release 3.0 mappable to 3,689 proteins with distinct gene names, with a median of 3 proteins per complex (S8 Fig). Thus, small protein complexes are widespread in the proteome and a large fraction of the proteome could be placed under complex post-transcriptional control that decouples protein levels from mRNA levels. Moreover, half of the proteins in the CORUM annotations belong to 2 or more distinct complexes (including intermediates and subcomplexes) which suggest further opportunities for more complex patterns of interdependent protein and mRNA levels.

Network representation of the interdependencies of protein level regulations

We next examined whether the interdependency of proteins and trans locus transcripts may be used to infer potential novel regulatory drivers. To do so, we constructed graphical models to visualize the connections between each protein and the transcriptome features that contribute to its predicted level, followed by topological analysis to find hub nodes and extract network patterns. To limit scope, we generated graphs to proteins whose prediction improved in the CPTAC_8 data set following the inclusion of the CORUM or STRING features ($\Delta r \geq 0.25$) A directed graph is then generated using a list of edges that connect feature variables

(transcripts) to target variables (proteins), with the edge weights calculated using a function of the elastic net coefficients and random forest feature importance of the models for both the CORUM (S6 Table) and the STRING (S7 Table) feature sets. The resulting network is partitioned into connected components of transcript-protein relationships. The subgraph that corresponds to the relationship between PCCA and PCCB as discussed above shows the PCCB node contributing (as transcript) to its own protein level weakly, and the PCCA node contributing strongly to the PCCB node (as protein) (Fig 6A). The relationship of CMC1 with other cytochrome c oxidase complex subunits, as described above, is represented in a wheel and spokes pattern where multiple transcripts (orange nodes) contribute either positively or negatively to the CMC1 protein level, with MT-CO1 having the strongest positive contributions (Fig 6B). Another member of the MITRAC module COA3 to which CMC1 binds likewise exerts a positive influence on the protein level of CMC1, consistent with the possibility that interdependent protein and transcript relationships may reflect information related to protein complex assembly sequences. Larger subgraphs represent instances where two or more protein target nodes are connected by shared transcript source nodes. Due to the method by which the graphs are generated, protein hubs (with high in-degrees) are more prevalent, but transcripts with high outflow activities are also seen that connect two or more protein nodes. A subgraph in Fig 6C predicts that the transcript levels of SH3KBP1 and CD2AP both contribute to the predicted protein levels of multiple proteins involved in the RICH1/AMOT polarity complex and the PI4K2A-WASH complex, and the protein level of WASHC1 is influenced by genes in the CCC-Wash (WASH1, FAM21C) complex and the PI4K2A-WASH complex (Fig 6C).

We highlight an instance where mRNA-protein relationships may be examined to generate new hypotheses. The mitochondrial ribosome is represented in a highly connected subgraph where the abundances of a majority of complex subunit proteins are each contingent upon multiple transcripts (Fig 6D). Using topological analysis based on the hierarchical community detection framework (HiDeF) method [33], we showed that this subgraph preserves the expected hierarchical relationship between the 28S small and 39S large subunits in forming the 55S mitochondrial ribosome, suggesting the preferential connections within the proteogenomic graphs broadly recapitulate subcomplex assembly (Fig 6E). We then extracted the hub nodes from the subgraph using the betweenness centrality algorithm, which revealed mitochondrial serine beta-lactamase-like protein (LACTB) to be an upstream gene that is highly connected within the subgraph and exerts a regulatory effect on both small and large mitochondrial ribosome subunit components. Feature importance analysis further suggests the LACTB transcript to have a negative impact on multiple ribosomal mitochondrial proteins (Fig 6F). Of note, LACTB was previously identified as a new subunit (MRP-L56) of the mitochondrial ribosomal large subunit isolated from sucrose gradients [58]. Subsequent work however has established the LACTB protein as a filament-forming protein that is localized to the intermembrane space instead away from the mitochondrial ribosomes [59] that possesses in vitro protease activity [60] and that acts as a tumor suppressor by maintaining post-mitotic differentiation states [60]. The physiological function and antiproliferative mechanism of LACTB are incompletely understood, although induced LACTB expression in cancer cells was found to affect mitochondrial phospholipid metabolism [60]. The analysis here therefore raises the possibility that LACTB may also affect oxidative phosphorylation in differentiated cells by regulating mitochondrial ribosome biogenesis and protein synthesis, which may be validated experimentally.

A similar hub node analysis from the largest subgraph from the STRING feature set (S8Table) showed chromogranin A and B (CHGA and CHGB) as inward hubs (i.e., proteins associated with many transcripts) associated with secreted peptides including CST3, SER-PINC1, MFGE8, SERPIND1, and others (S9 Fig), which is consistent with the known roles of

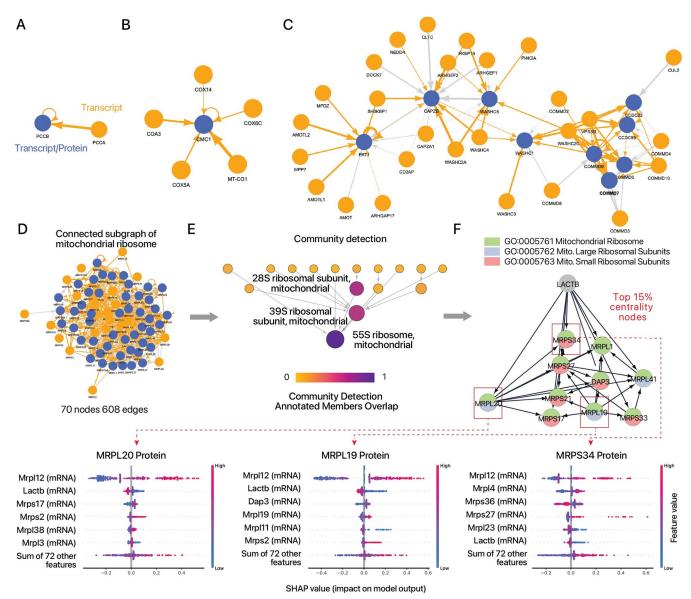


Fig 6. Directed graphs of protein and transcript interrelationships identify candidate regulatory genes. A-C. Examples of directed graphs constructed from genome-wide relationships of transcript-predicted proteins, containing members of A. the propionyl-CoA carboxylase complex; B. the cytochrome c oxidase, mitochondrial complex; C. the PI4K2A-WASH complex, the RICH1/AMOT polarity complex, and others. In each subgraph, orange nodes have outflow edges only (i.e., they are contributing transcripts in the prediction models). Blue nodes are nodes that are connected to other nodes via at least one inflow edge (i.e., they represent proteins, and optionally also transcripts if they also have outward edges). Orange edges represent positive coefficients of the transcripts to the target proteins in the elastic net models; gray edges represent negative coefficients. All edges are directed from transcript to protein, and the widths of the edges are scaled by the weight. D. A highly connected subgraph of mitochondrial ribosome subunits containing 73 nodes and 834 edges. E. Persistent community detection and network representation of preferential node connections, showing a hierarchical relationship between the 28S and 39S subcomplex with the assembled 55S mitochondrial ribosome. F. Network representation of hub nodes defined as 15% of nodes ranked by betweenness centrality, which predicts a potential role of LACTB as a critical hub that lies upstream of multiple large and small mitochondrial ribosomal protein subunits. Node colors represent the pie chart diagram of the corresponding GO biological process described in the table. SHAP values of three proteins (MRPL20, MRPL19, MRPS34) are highlighted showing top model contributors.

https://doi.org/10.1371/journal.pcbi.1010702.g006

these acidic glycoproteins as the primary constituents of secretory granules that can regulate their rate of their formation [61]. Taken together, the topological analysis shown here suggests mRNA-protein networks may be useful for generating new hypothesis on regulatory drivers from large proteogenomics data.

Discussion

Widespread contributions of trans locus transcripts to protein level

Recent work has applied various machine learning models to the computational task of predicting across-sample protein levels using transcriptomics data, but investigation into the biological factors that uncouple transcriptome and proteome data have remained limited, with many studies noting only that predictability may differ across broad functional categories, e.g., essential genes or metabolic genes may be less predictable. Here we analyzed cases where proteins are poorly predicted by their cognate mRNA, and showed that the transcript levels of interacting partners had an outsized contribution to the abundance of proteins of interest. Although the notion that protein-protein interaction can influence protein abundance posttranscriptionally is not new, as it is known that supernumerary subunits of protein complexes can be removed through protein degradation, details of genes that show especially poor mRNA and protein abundance correlations and the identities of their protein regulators have remained scarce. The current study provides new evidence from protein prediction models that protein abundance by other transcripts is common in the proteome, nominating specific interactions involving not only large megadalton sized multi-protein complexes as previously observed, but also smaller stable complexes (e.g., propionyl-CoA carboxylase with two subunits PCCB and PCCA; calcineurin A and B) in the CORUM feature set and potentially more transient interactions documented in the STRING feature set. The common presence of small stable complexes in the proteome greatly expands the repertoire of proteins for which due considerations should be given when directly interpreting transcript level data as representing protein level information. Moreover, it has been suggested that promiscuous protein-protein interactions without established biological function may be a common occurrence during the co-evolution of functional protein-protein interactions [62]. The occurrence of such stable but non-functional interacting pairs could further increase the scope of trans locus protein regulation which would have implications on the predictability of protein levels from mRNA levels.

The imperfect correlation between protein and mRNA points to orthogonal information that exists between transcript and protein regulations, which underpins the untapped potential for further multi-omics integration to derive new insights [63]. The use of protein correlations to find causal insights for post-transcriptional regulations has previously been explored in other contexts, which looked for anti-correlation between E3 ubiquitin ligases with known proteins or ubiquitination sites of interest, which may control their protein level by virtue of post-translational degradation [18,24]. Extending this intuition towards multi-omics correlations, we used directed graphs generated from the model coefficients and feature importance of trans locus transcripts to represent the interaction patterns between the transcripts and proteins of interest. This in turn nominated a number of hub proteins whose abundance is contributed by multiple transcripts and hub transcripts that regulate multiple proteins, and enabled community detection analysis to find the relationships between biological processes in protein and mRNA correlation networks. As more data continue to become available, we foresee that graphical models will be useful for finding more trans locus regulations of protein levels, such as those that represent known assembly sequential steps or post-transcriptional regulators. These graphical models may be used to generate testable hypotheses or find utility in predicting experimental outcome, barring confounders or reverse causality effects. For instance, in the case of the actin related protein 2/3 complex, one might predict that within a certain concentration range an overexpression of ACTR3 and ARPC4 will be more effective in modulating ARPC3 protein levels than augmenting the expression of ARPC3 itself, which is readily testable by experimentation. More generally, the results here corroborate the importance of post-transcriptional regulation including protein degradation and turnover in modulating protein levels.

Limitations of the study

Significant differences in protein regulations likely occur in different cell and tissue types. The predictive models here are trained using publicly available CPTAC data from 8 cancer types, which contain transcriptome and proteome data from both tumor and normal adjacent tissue samples. In prior work, we found general concordance between protein and mRNA correlation in CPTAC samples vs. GTEx tissue proteomics data [15]. However, other studies that compared protein and mRNA correlation in tumors and normal adjacent tissues have found higher inter-sample correlation in tumor samples [21,22], which may be attributable to the increased translation rates in cancer. Hence, additional discordant cases between proteins and mRNA likely remain to be discovered that are omitted here. Sample difference may also explain the observation from GTEx proteomics that secreted proteins are associated with protein and mRNA discordance [7] as one might intuitively assume but which is not apparent in the analysis here. The cases used for training the model are not labeled by their cancer type or their tumor vs. normal adjacent tissue designation, which would likely have availed overall predictive performance.

We performed the analysis using expression data from TCGA/CPTAC RNA sequencing experiments through the cptac Python API, which retrieves the final data tables from the flagship CPTAC papers of each individual cancer type [25]. Although each TCGA/CPTAC cancer subtype project follows an overall consistent experimental design and data acquisition strategy, minute differences exist in the processing pipelines used to analyze the RNA sequencing data (e.g., STAR vs. Bowtie2) and gene expression measure (e.g., RPKM vs. FPKM) which could bias gene expression values across cancer types. Likewise, there are subtle differences in the protein expression data in each cancer type (e.g., search engine, isotope tags). As our goal here is not to compare trends across cancer types, we have taken inspiration from the CPTAC Dream Challenge submissions that improved overall predictions by borrowing information across cancer types. However, it is also known that different data sets may present different variability. Future work may employ uniformed processed data to improve performance and reduce bias, such as data from the TCGA PanCanAtlas which processed all TCGA cancer RNA sequencing data uniformly [64], or compare batch correction strategies. To our knowledge, uniformly processed protein data are not yet available at the time of writing. Other machine learning and deep learning algorithms can likely further boost protein level predictions. We have limited our scope here to comparisons of different feature sets and previously employed algorithms, and to result interpretation.

Finally, pitfalls must be considered when attempting to interpret predictive models in search of mechanistic insights. Feature importance in the predictive models represents correlation rather than causality, and hence interpretations can be confounded by independent confounders, multicollinearity, and reverse causality. In simple cases such as PCCB–PCCA, given the emphasis on prior feature selections that prioritize known protein-protein interaction partners, the logical interpretation would be to assume the hierarchical nature of gene regulations where transcript levels are more likely to affect protein abundance than the opposite, but this assumption becomes more tenuous as feature sets expand and the number of associations increase, and in cases such as transcription factor proteins whose abundance can affect transcript levels across samples. Future work in this area may employ more sophisticated causal inference methods to identify regulatory modalities.

Conclusion

In summary, this study compared predictive models of protein levels using different transcript feature sets, and provided biological interpretations of the results by highlighting trans locus

transcripts with substantial contributions to protein levels. The analysis here therefore reveals new details into the gene identity and modality of trans regulation of protein levels, and gives support to further development of prior transcript feature selection strategies to optimize protein prediction tasks. The results show that the transcript levels of protein-protein interaction partners can broadly influence protein abundance in a tissue, which has implications on the interpretations of transcriptomics data and on understanding the architecture of proteome composition regulations. With further refinement of feature selection and feature engineering methods and the availability of large data sets, we foresee that similar approaches to those shown here will provide valuable new insights into post-transcriptional mechanisms of protein regulations.

Supporting information

S1 Fig. Model performance measured by additional metrics. A. Box plots of test set normalized root mean square error (NRMSE) between the transcript-predicted and actual protein level for each protein are shown across five feature sets (column: single/self transcript, CORUM interactors, STRING high-confidence associated proteins; STRING low-confidence associated proteins, and all transcripts) and three algorithms (multiple linear regression, elastic net, and random forest). In each plot, x axis denotes the number of CPTAC data set used to train the models box: interquartile range; whiskers: +/- 1.5 IQR. **B.** As above, but for test set goodness-of-fit (R²). (PDF)

S2 Fig. A. Model performance in single cancer data sets. Box plots of test set correlation coefficients (r) between the transcript-predicted and actual protein level for each protein are shown across five feature sets (column: single/self transcript, CORUM interactors, STRING high-confidence associated proteins; STRING low-confidence associated proteins, and all transcripts) and three algorithms (multiple linear regression, elastic net, and random forest). In each plot, x axis denotes the CPTAC cancer type study used to train the models; box: interquartile range; whiskers: +/- 1.5 IQR. **B.** Model performance when the 8 data sets were combined in the order of decreasing single data set performance. (PDF)

S3 Fig. Improvements to protein prediction from the incorporation of additional transcript features. A-B. Scatterplot showing a significant linear relationship between the number of low-confidence interactors a protein has as annotated in STRING vs. the increase in test set correlation between predicted vs. actual protein levels in the Elastic Net models over the self-transcript feature set (correlation test p: 8.2e–131 for all proteins; 6.5e–130 for proteins with 10 or more interactors). **C-D.** As above, but for the CORUM feature set. (PDF)

S4 Fig. Additional proteins with improved predicted levels after inclusion of additional transcript features. Two additional highlighted proteins with substantial predictability from transcriptome data upon the inclusion of additional features are shown: **A.** MICU2, **B.** PPP3R1. For each protein, the transcript-trained prediction of protein level is plotted on the x axis and the actual protein level is plotted on the y axis. Blue: train set, brown: test set. Columns denote the transcript feature set used to train the model. The number of features used to train the model in each feature set is shown inside each plot. r: Correlation coefficient. (PDF)

S5 Fig. Cross-omics co-expression of MICU2 and PPP3R1 with functionally associated proteins. Two examples of proteins whose abundance is better explained by another transcript

are shown. **A.** MICU2 protein level is predicted by MCU transcript but not its own transcript. **B.** PPP3R1 protein level is explained by PPP3CA transcript level but not its own transcript. Substantial correlations across transcripts and proteins (\geq 0.4) are bolded. (PDF)

S6 Fig. SHAP interpretation of feature importance in **A.** random forest and **B.** gradient boosting model output from the CPTAC_8 CORUM feature set. The SHAP values of top transcript features and their impact on model output are shown for each of six proteins highlighted in the text (from top to bottom): PCCB, CMC1, PSMG2, SMCR8, MICU2, PPP3R1. (PDF)

S7 Fig. SHAP interpretation of feature importance in the CORUM feature set, random forest model, of a single data set (LSCC), showing largely conserved observations among the highlighted proteins PCCB, CMC1, PSMG2, SMCR8, MICU2, PPP3R1. Consistent with the analysis of the combined CPTAC_8 data set, five of the six proteins are not best predicted by their cognate transcripts and all but one (MT-CO1) of the top trans locus transcripts are preserved, which was not among the examined candidate features. Only test set data points are shown. (PDF)

S8 Fig. Density plot showing the distribution of complex sizes (x:-axis: number of members in complex; y-axis: log10 total polypeptide molecular weight in complex) in the annotated feature set derived from CORUM v.3.0. Names of select complexes are labeled. The majority of complexes are small with a median of 3 proteins per complex. (PDF)

S9 Fig. The most connected subgraph in the STRING feature set contains 6,319 nodes and 11,161 edges, containing proteins belonging to multiple distinct cellular compartments and multi-protein complexes. The hierarchical network diagram shows Chromogranin A and B (CHGA/CHGB) (magenta) and their first-degree inward flow neighbors (green) in the largest subgraph constructed from the STRING data set. The protein level of CHGA and CHGB is associated with the transcript levels of multiple secreted peptide coding genes, consistent with these proteins serve to regulate gene expression in secretion pathways. (PDF)

S1 Table. Table listing performance metrics (test set correlation coefficients, R², NRMSE) of the tested models, feature sets, and data sets.

(TSV)

S2 Table. Table listing Gene Ontology cellular compartment term enrichment results for proteins that are well predicted by their transcript levels.

(TSV)

S3 Table. Table listing Gene Ontology cellular compartment term enrichment for proteins that are poorly predicted by their transcript levels.
(TSV)

S4 Table. Table listing individual genewise model performance metrics (test set correlation coefficients, R², NRMSE) in the CPTAC_8 data set using elastic net models on different feature sets.

(TSV)

S5 Table. Table listing Gene Ontology cellular compartment term enrichment for proteins with substantial improvements in prediction performance upon incorporating STRING

features in elastic net models.

(TSV)

S6 Table. Network edge list for proteins with improved predictions in the CORUM feature set along with the coefficients and feature importances of contributing transcripts. (CSV)

S7 Table. Network edge list for proteins with improved predictions in the STRING feature set along with the coefficients and feature importances of contributing transcripts. (CSV)

S8 Table. Hub genes based on centrality scores in the STRING feature set network. (TXT)

Author Contributions

Conceptualization: Himangi Srivastava, Jordan Currie, Maggie P. Y. Lam, Edward Lau.

Formal analysis: Himangi Srivastava, Michael J. Lippincott, Edward Lau.

Funding acquisition: Maggie P. Y. Lam, Edward Lau.

Investigation: Michael J. Lippincott, Robert Canfield, Maggie P. Y. Lam, Edward Lau.

Resources: Maggie P. Y. Lam, Edward Lau. Software: Himangi Srivastava, Edward Lau. Supervision: Maggie P. Y. Lam, Edward Lau.

Writing - original draft: Himangi Srivastava, Jordan Currie, Maggie P. Y. Lam, Edward Lau.

Writing – review & editing: Michael J. Lippincott, Jordan Currie, Maggie P. Y. Lam, Edward Lau.

References

- Gygi SP, Rochon Y, Franza BR, Aebersold R. Correlation between protein and mRNA abundance in yeast. Mol Cell Biol. 1999; 19: 1720–1730. https://doi.org/10.1128/MCB.19.3.1720 PMID: 10022859
- Liu Y, Beyer A, Aebersold R. On the Dependency of Cellular Protein Levels on mRNA Abundance. Cell. 2016; 165: 535–550. https://doi.org/10.1016/j.cell.2016.03.014 PMID: 27104977
- Vogel C, Marcotte EM. Insights into the regulation of protein abundance from proteomic and transcriptomic analyses. Nat Rev Genet. 2012; 13: 227–232. https://doi.org/10.1038/nrg3185 PMID: 22411467
- Franks A, Airoldi E, Slavov N. Post-transcriptional regulation across human tissues. PLoS Comput Biol. 2017; 13: e1005535. https://doi.org/10.1371/journal.pcbi.1005535 PMID: 28481885
- Upadhya SR, Ryan CJ. Experimental reproducibility limits the correlation between mRNA and protein abundances in tumour proteomic profiles. Systems Biology; 2021 Sep. https://doi.org/10.1101/2021.09.22.461108
- An H, Harper JW. Ribosome Abundance Control Via the Ubiquitin–Proteasome System and Autophagy. J Mol Biol. 2020; 432: 170–184. https://doi.org/10.1016/j.jmb.2019.06.001 PMID: 31195016
- Jiang L, Wang M, Lin S, Jian R, Li X, Chan J, et al. A Quantitative Proteome Map of the Human Body. Cell. 2020; 183: 269–283.e19. https://doi.org/10.1016/j.cell.2020.08.036 PMID: 32916130
- Gonçalves E, Fragoulis A, Garcia-Alonso L, Cramer T, Saez-Rodriguez J, Beltrao P. Widespread Post-transcriptional Attenuation of Genomic Copy-Number Variation in Cancer. Cell Syst. 2017; 5: 386–398.
 https://doi.org/10.1016/j.cels.2017.08.013 PMID: 29032074
- Eraslan B, Wang D, Gusic M, Prokisch H, Hallström BM, Uhlén M, et al. Quantification and discovery of sequence determinants of protein-per-mRNA amount in 29 human tissues. Mol Syst Biol. 2019; 15: e8513. https://doi.org/10.15252/msb.20188513 PMID: 30777893

- Wang D, Eraslan B, Wieland T, Hallström B, Hopf T, Zolg DP, et al. A deep proteome and transcriptome abundance atlas of 29 healthy human tissues. Mol Syst Biol. 2019;15. https://doi.org/10.15252/msb.20188503 PMID: 30777892
- Waldbauer JR, Rodrigue S, Coleman ML, Chisholm SW. Transcriptome and proteome dynamics of a light-dark synchronized bacterial cell cycle. PloS One. 2012; 7: e43432. https://doi.org/10.1371/journal.pone.0043432 PMID: 22952681
- Yang M, Petralia F, Li Z, Li H, Ma W, Song X, et al. Community Assessment of the Predictability of Cancer Protein and Phosphoprotein Levels from Genomics and Transcriptomics. Cell Syst. 2020; 11: 186–195.e9. https://doi.org/10.1016/j.cels.2020.06.013 PMID: 32710834
- Eicher T, Patt A, Kautto E, Machiraju R, Mathé E, Zhang Y. Challenges in proteogenomics: a comparison of analysis methods with the case study of the DREAM proteogenomics sub-challenge. BMC Bioinformatics. 2019; 20: 669. https://doi.org/10.1186/s12859-019-3253-z PMID: 31861998
- Li H, Siddiqui O, Zhang H, Guan Y. Joint learning improves protein abundance prediction in cancers. BMC Biol. 2019; 17: 107. https://doi.org/10.1186/s12915-019-0730-9 PMID: 31870366
- Han Y, Li LZ, Kastury NL, Thomas CT, Lam MPY, Lau E. Transcriptome features of striated muscle aging and predictability of protein level changes. Mol Omics. 2021. https://doi.org/10.1039/d1mo00178g PMID: 34328155
- Xu W, He H, Guo Z, Li W. Evaluation of machine learning models on protein level inference from prioritized RNA features. Brief Bioinform. 2022; bbac091. https://doi.org/10.1093/bib/bbac091 PMID: 35352096
- Zhang H, Liu T, Zhang Z, Payne SH, Zhang B, McDermott JE, et al. Integrated Proteogenomic Characterization of Human High-Grade Serous Ovarian Cancer. Cell. 2016; 166: 755–765. https://doi.org/10.1016/j.cell.2016.05.069 PMID: 27372738
- Mertins P, Mani DR, Ruggles KV, Gillette MA, Clauser KR, Wang P, et al. Proteogenomics connects somatic mutations to signalling in breast cancer. Nature. 2016; 534: 55–62. https://doi.org/10.1038/nature18003 PMID: 27251275
- Dou Y, Kawaler EA, Cui Zhou D, Gritsenko MA, Huang C, Blumenberg L, et al. Proteogenomic Characterization of Endometrial Carcinoma. Cell. 2020; 180: 729–748.e26. https://doi.org/10.1016/j.cell.2020.01.026 PMID: 32059776
- Vasaikar S, Huang C, Wang X, Petyuk VA, Savage SR, Wen B, et al. Proteogenomic Analysis of Human Colon Cancer Reveals New Therapeutic Opportunities. Cell. 2019; 177: 1035–1049.e19. https://doi.org/10.1016/j.cell.2019.03.030 PMID: 31031003
- Gillette MA, Satpathy S, Cao S, Dhanasekaran SM, Vasaikar SV, Krug K, et al. Proteogenomic Characterization Reveals Therapeutic Vulnerabilities in Lung Adenocarcinoma. Cell. 2020; 182: 200–225.e35. https://doi.org/10.1016/j.cell.2020.06.013 PMID: 32649874
- Clark DJ, Dhanasekaran SM, Petralia F, Pan J, Song X, Hu Y, et al. Integrated Proteogenomic Characterization of Clear Cell Renal Cell Carcinoma. Cell. 2019; 179: 964–983.e31. https://doi.org/10.1016/j.cell.2019.10.007 PMID: 31675502
- 23. Wang L-B, Karpova A, Gritsenko MA, Kyle JE, Cao S, Li Y, et al. Proteogenomic and metabolomic characterization of human glioblastoma. Cancer Cell. 2021; 39: 509–528.e20. https://doi.org/10.1016/j.ccell.2021.01.006 PMID: 33577785
- Satpathy S, Krug K, Jean Beltran PM, Savage SR, Petralia F, Kumar-Sinha C, et al. A proteogenomic portrait of lung squamous cell carcinoma. Cell. 2021; 184: 4348–4371.e40. https://doi.org/10.1016/j.cell.2021.07.016 PMID: 34358469
- Lindgren CM, Adams DW, Kimball B, Boekweg H, Tayler S, Pugh SL, et al. Simplified and Unified Access to Cancer Proteogenomic Data. J Proteome Res. 2021; 20: 1902–1910. https://doi.org/10.1021/acs.jproteome.0c00919 PMID: 33560848
- Giurgiu M, Reinhard J, Brauner B, Dunger-Kaltenbach I, Fobo G, Frishman G, et al. CORUM: the comprehensive resource of mammalian protein complexes-2019. Nucleic Acids Res. 2019; 47: D559

 D563. https://doi.org/10.1093/nar/gky973 PMID: 30357367
- Szklarczyk D, Gable AL, Nastou KC, Lyon D, Kirsch R, Pyysalo S, et al. The STRING database in 2021: customizable protein-protein networks, and functional characterization of user-uploaded gene/measurement sets. Nucleic Acids Res. 2021; 49: D605–D612. https://doi.org/10.1093/nar/gkaa1074 PMID: 33237311
- 28. Pedregosa F, Varoquaux G, Gramfort A, Michel V, Thirion B, Grisel O, et al. Scikit-learn: Machine Learning in Python. J Mach Learn Res. 2011; 12: 2825–2830.
- 29. Lundberg SM, Lee S-I. A Unified Approach to Interpreting Model Predictions. In: Guyon I, Luxburg UV, Bengio S, Wallach H, Fergus R, Vishwanathan S, et al., editors. Advances in Neural Information

- Processing Systems. Curran Associates, Inc.; 2017. Available: https://proceedings.neurips.cc/paper/2017/file/8a20a8621978632d76c43dfd28b67767-Paper.pdf
- Kursa MB, Rudnicki WR. Feature Selection with the Boruta Package. J Stat Softw. 2010;36. https://doi. org/10.18637/jss.v036.i11
- Hagberg A, Swart P, S Chult D. Exploring network structure, dynamics, and function using networkx. Los Alamos National Lab. (LANL), Los Alamos, NM (United States); 2008 Jan. Report No.: LA-UR-08-05495; LA-UR-08-5495. Available: https://www.osti.gov/biblio/960616-exploring-network-structure-dynamics-function-using-networkx
- Shannon P, Markiel A, Ozier O, Baliga NS, Wang JT, Ramage D, et al. Cytoscape: a software environment for integrated models of biomolecular interaction networks. Genome Res. 2003; 13: 2498–2504. https://doi.org/10.1101/gr.1239303 PMID: 14597658
- Zheng F, Zhang S, Churas C, Pratt D, Bahar I, Ideker T. HiDeF: identifying persistent structures in multiscale 'omics data. Genome Biol. 2021; 22: 21. https://doi.org/10.1186/s13059-020-02228-4 PMID: 33413539
- Singhal A, Cao S, Churas C, Pratt D, Fortunato S, Zheng F, et al. Multiscale community detection in Cytoscape. Przytycka TM, editor. PLOS Comput Biol. 2020; 16: e1008239. https://doi.org/10.1371/journal.pcbi.1008239 PMID: 33095781
- 35. Traag VA, Waltman L, van Eck NJ. From Louvain to Leiden: guaranteeing well-connected communities. Sci Rep. 2019; 9: 5233. https://doi.org/10.1038/s41598-019-41695-z PMID: 30914743
- Reimand J, Arak T, Adler P, Kolberg L, Reisberg S, Peterson H, et al. g:Profiler—a web server for functional interpretation of gene lists (2016 update). Nucleic Acids Res. 2016; 44: W83–W89. https://doi.org/10.1093/nar/gkw199 PMID: 27098042
- Tang Y, Li M, Wang J, Pan Y, Wu F-X. CytoNCA: A cytoscape plugin for centrality analysis and evaluation of protein interaction networks. Biosystems. 2015; 127: 67–72. https://doi.org/10.1016/j.biosystems.2014.11.005 PMID: 25451770
- Doncheva NT, Morris JH, Gorodkin J, Jensen LJ. Cytoscape StringApp: Network Analysis and Visualization of Proteomics Data. J Proteome Res. 2019; 18: 623–632. https://doi.org/10.1021/acs.jproteome. 8b00702 PMID: 30450911
- **39.** Yu G, He Q-Y. ReactomePA: an R/Bioconductor package for reactome pathway analysis and visualization. Mol Biosyst. 2016; 12: 477–479. https://doi.org/10.1039/c5mb00663e PMID: 26661513
- Wu T, Hu E, Xu S, Chen M, Guo P, Dai Z, et al. clusterProfiler 4.0: A universal enrichment tool for interpreting omics data. The Innovation. 2021; 2: 100141. https://doi.org/10.1016/j.xinn.2021.100141 PMID: 34557778
- 41. Gu Z, Gu L, Eils R, Schlesner M, Brors B. circlize implements and enhances circular visualization in R. Bioinformatics. 2014; 30: 2811–2812. https://doi.org/10.1093/bioinformatics/btu393 PMID: 24930139
- Waskom M. seaborn: statistical data visualization. J Open Source Softw. 2021; 6: 3021. https://doi.org/10.21105/joss.03021
- Mick DU, Dennerlein S, Wiese H, Reinhold R, Pacheu-Grau D, Lorenzi I, et al. MITRAC Links Mitochondrial Protein Translocation to Respiratory-Chain Assembly and Translational Regulation. Cell. 2012; 151: 1528–1541. https://doi.org/10.1016/j.cell.2012.11.053 PMID: 23260140
- Signes A, Fernandez-Vizarra E. Assembly of mammalian oxidative phosphorylation complexes I–V and supercomplexes. Essays Biochem. 2018; 62: 255–270. https://doi.org/10.1042/EBC20170098 PMID: 30030361
- 45. Hirano Y, Hendil KB, Yashiroda H, Iemura S, Nagane R, Hioki Y, et al. A heterodimeric complex that promotes the assembly of mammalian 20S proteasomes. Nature. 2005; 437: 1381–1385. https://doi.org/10.1038/nature04106 PMID: 16251969
- 46. The GTEx Consortium. The GTEx Consortium atlas of genetic regulatory effects across human tissues. Science. 2020; 369: 1318–1330. https://doi.org/10.1126/science.aaz1776 PMID: 32913098
- 47. Sasaki K, Hamazaki J, Koike M, Hirano Y, Komatsu M, Uchiyama Y, et al. PAC1 gene knockout reveals an essential role of chaperone-mediated 20S proteasome biogenesis and latent 20S proteasomes in cellular homeostasis. Mol Cell Biol. 2010; 30: 3864–3874. https://doi.org/10.1128/MCB.00216-10 PMID: 20498273
- 48. Tang D, Sheng J, Xu L, Zhan X, Liu J, Jiang H, et al. Cryo-EM structure of C9ORF72–SMCR8–WDR41 reveals the role as a GAP for Rab8a and Rab11a. Proc Natl Acad Sci. 2020; 117: 9876–9883. https://doi.org/10.1073/pnas.2002110117 PMID: 32303654
- Sellier C, Campanari M-L, Julie Corbier C, Gaucherot A, Kolb-Cheynel I, Oulad-Abdelghani M, et al. Loss of C9ORF72 impairs autophagy and synergizes with polyQ Ataxin-2 to induce motor neuron dysfunction and cell death. EMBO J. 2016; 35: 1276–1297. https://doi.org/10.15252/embj.201593350 PMID: 27103069

- 50. Sullivan PM, Zhou X, Robins AM, Paushter DH, Kim D, Smolka MB, et al. The ALS/FTLD associated protein C9orf72 associates with SMCR8 and WDR41 to regulate the autophagy-lysosome pathway. Acta Neuropathol Commun. 2016; 4: 51. https://doi.org/10.1186/s40478-016-0324-5 PMID: 27193190
- Kamer KJ, Mootha VK. MICU1 and MICU2 play nonredundant roles in the regulation of the mitochondrial calcium uniporter. EMBO Rep. 2014; 15: 299–307. https://doi.org/10.1002/embr.201337946 PMID: 24503055
- Plovanich M, Bogorad RL, Sancak Y, Kamer KJ, Strittmatter L, Li AA, et al. MICU2, a paralog of MICU1, resides within the mitochondrial uniporter complex to regulate calcium handling. PloS One. 2013; 8: e55785. https://doi.org/10.1371/journal.pone.0055785 PMID: 23409044
- Hoffman A, Taleski G, Sontag E. The protein serine/threonine phosphatases PP2A, PP1 and calcineurin: A triple threat in the regulation of the neuronal cytoskeleton. Mol Cell Neurosci. 2017; 84: 119–131. https://doi.org/10.1016/j.mcn.2017.01.005 PMID: 28126489
- Rusnak F, Mertz P. Calcineurin: form and function. Physiol Rev. 2000; 80: 1483–1521. https://doi.org/ 10.1152/physrev.2000.80.4.1483 PMID: 11015619
- Liberzon A, Birger C, Thorvaldsdóttir H, Ghandi M, Mesirov JP, Tamayo P. The Molecular Signatures Database (MSigDB) hallmark gene set collection. Cell Syst. 2015; 1: 417–425. https://doi.org/10.1016/j.cels.2015.12.004 PMID: 26771021
- 56. Lau E, Venkatraman V, Thomas CT, Wu JC, Van Eyk JE, Lam MPY. Identifying High-Priority Proteins Across the Human Diseasome Using Semantic Similarity. J Proteome Res. 2018; 17: 4267–4278. https://doi.org/10.1021/acs.jproteome.8b00393 PMID: 30256117
- 57. Lam MPY, Venkatraman V, Xing Y, Lau E, Cao Q, Ng DCM, et al. Data-Driven Approach To Determine Popular Proteins for Targeted Proteomics Translation of Six Organ Systems. J Proteome Res. 2016; 15: 4126–4134. https://doi.org/10.1021/acs.jproteome.6b00095 PMID: 27356587
- Koc EC, Burkhart W, Blackburn K, Moyer MB, Schlatzer DM, Moseley A, et al. The Large Subunit of the Mammalian Mitochondrial Ribosome: ANALYSIS OF THE COMPLEMENT OF RIBOSOMAL PRO-TEINS PRESENT*. J Biol Chem. 2001; 276: 43958–43969. https://doi.org/10.1074/jbc.M106510200 PMID: 11551941
- 59. Polianskyte Z, Peitsaro N, Dapkunas A, Liobikas J, Soliymani R, Lalowski M, et al. LACTB is a filament-forming protein localized in mitochondria. Proc Natl Acad Sci. 2009; 106: 18960–18965. https://doi.org/10.1073/pnas.0906734106 PMID: 19858488
- 60. Keckesova Z, Donaher JL, De Cock J, Freinkman E, Lingrell S, Bachovchin DA, et al. LACTB is a tumour suppressor that modulates lipid metabolism and cell state. Nature. 2017; 543: 681–686. https:// doi.org/10.1038/nature21408 PMID: 28329758
- Kim T, Tao-Cheng J-H, Eiden LE, Loh YP. Chromogranin A, an "On/Off" Switch Controlling Dense-Core Secretory Granule Biogenesis. Cell. 2001; 106: 499–509. https://doi.org/10.1016/s0092-8674(01) 00459-7 PMID: 11525735
- Aakre CD, Herrou J, Phung TN, Perchuk BS, Crosson S, Laub MT. Evolving new protein-protein interaction specificity through promiscuous intermediates. Cell. 2015; 163: 594–606. https://doi.org/10.1016/j.cell.2015.09.055 PMID: 26478181
- Payne SH. The utility of protein and mRNA correlation. Trends Biochem Sci. 2015; 40: 1–3. https://doi. org/10.1016/j.tibs.2014.10.010 PMID: 25467744
- Sanchez-Vega F, Mina M, Armenia J, Chatila WK, Luna A, La KC, et al. Oncogenic Signaling Pathways in The Cancer Genome Atlas. Cell. 2018; 173: 321–337.e10. https://doi.org/10.1016/j.cell.2018.03.035 PMID: 29625050