



Published in final edited form as:

Mol Cell. 2016 October 20; 64(2): 282–293. doi:10.1016/j.molcel.2016.09.003.

SONAR discovers RNA binding proteins from analysis of large-scale protein-protein interactomes

Kristopher W. Brannan^{1,§}, Wenhao Jin^{2,§}, Stephanie C. Huelga^{1,§}, Charles A. S. Banks³, Joshua M. Gilmore³, Laurence Florens³, Michael P. Washburn^{3,4}, Eric L. Van Nostrand¹, Gabriel A. Pratt¹, Marie K. Schwinn⁵, Danette L. Daniels⁵, and Gene W. Yeo^{1,2,6}

¹Department of Cellular and Molecular Medicine, Stem Cell Program and Institute for Genomic Medicine; University of California, San Diego; La Jolla, California, 92093; USA

²Department of Physiology, Yong Loo Lin School of Medicine, National University of Singapore, Singapore

³Stowers Institute for Medical Research, 1000 E. 50th Street, Kansas City, MO 64110, USA

⁴Department of Pathology and Laboratory Medicine, University of Kansas Medical Center, 3901 Rainbow Boulevard, Kansas City, KS 66160, USA

⁵Promega Corporation, 2800 Woods Hollow Road, Madison, WI 53711 U.S.A

⁶Molecular Engineering Laboratory, A*STAR, Singapore

SUMMARY

RNA metabolism is controlled by an expanding yet incomplete catalog of RNA binding proteins (RBPs), many of which lack characterized RNA binding domains. Approaches to expand the RBP repertoire to discover non-canonical RBPs are currently needed. Here, HaloTag fusion pull-down of twelve nuclear and cytoplasmic RBPs followed by quantitative mass-spectrometry (MS) demonstrates that proteins interacting with multiple RBPs in an RNA-dependent manner are enriched for RBPs. This motivated SONAR, a computational approach that predicts RNA binding activity by analyzing large-scale affinity precipitation-MS protein-protein interactomes. Without relying on sequence or structure information, SONAR identifies 1923 human, 489 fly and 745

Correspondence: geneyeo@ucsd.edu.

[§]These authors contributed equally.

Lead Contact: Gene Yeo, geneyeo@ucsd.edu

Publisher's Disclaimer: This is a PDF file of an unedited manuscript that has been accepted for publication. As a service to our customers we are providing this early version of the manuscript. The manuscript will undergo copyediting, typesetting, and review of the resulting proof before it is published in its final citable form. Please note that during the production process errors may be discovered which could affect the content, and all legal disclaimers that apply to the journal pertain.

AUTHOR CONTRIBUTIONS

Conceptualization, K.W.B., W.J., S.C.H., D.L.D. and G.W.Y.; Investigation, K.W.B., W.J., S.C.H., D.L.D., M.P.W., C.A.S.B., J.M.G., L.F., and M.K.S.; Validation, K.W.B., D.L.D., M.P.W.; Formal Analysis, K.W.B., W.J., S.C.H., E.L.V.N., M.P.W., and G.A.P.; Writing – Original Draft, K.W.B., W.J., S.C.H., and G.W.Y.; Writing – Review & Editing, K.W.B., W.J., S.C.H., D.L.D., M.P.W., and G.W.Y.; Funding Acquisition G.W.Y.; Supervision G.W.Y.

ACCESSION NUMBERS

The accession number for candidate RBP eCLIPs reported in this paper is GEO: GSE86035

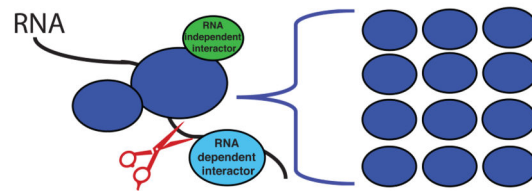
SUPPLEMENTAL INFORMATION

Supplemental Information includes Supplemental Experimental Procedures, six figures, and six tables and can be found with this article online at <http://xxx>

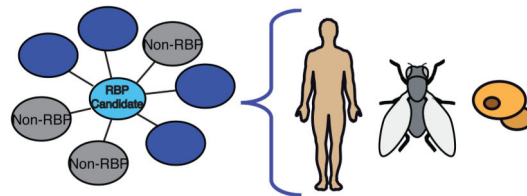
yeast RBPs, including over 100 human candidate RBPs that contain zinc finger domains. Enhanced CLIP confirms RNA binding activity and identifies transcriptome-wide RNA binding sites for SONAR-predicted RBPs, revealing unexpected RNA binding activity for disease-relevant proteins and DNA binding proteins.

Graphical Abstract

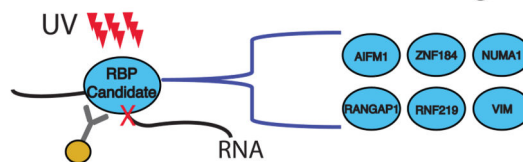
RBP interactomes are enriched for RBPs.



SONAR RBP predicts RBPs across species.



eCLIP identifies candidate RBP targets.



INTRODUCTION

RNA binding proteins (RBPs) act in both the nucleus and cytoplasm during every step of the RNA life cycle to exert precise and responsive control of gene expression. Given the number and expression levels of eukaryotic RBPs, and the fact that many are synthetically or embryonic lethal, it is unsurprising that evidence demonstrating their importance in development and disease emerges daily (Castello et al., 2013; Lukong et al., 2008).

The catalog of RBPs is diverse, and has grown to include metabolic enzymes, cell cycle regulators, and many other factors not previously associated with RNA biology (Beckmann et al., 2015; Conrad et al., 2016). As this diverse RBP catalog grows, complex regulatory networks are emerging that rely on highly coordinated and concurrent target association by numerous RBPs. With the advent of high-throughput techniques to accurately and robustly identify RBP target binding sites (Van Nostrand et al., 2016) and protein interaction partners, it has become possible to elucidate the complete list and the precise molecular functions of RBPs in these emerging RBP-RNA networks in a scalable manner.

Initial efforts to define the complete repertoire of RBPs largely focused on proteins with known or predicted RNA-binding domains (RBDs). Attempts to predict RBPs from primary

sequence and protein structure have yielded only moderate success (Si et al., 2015), particularly since more than a third of RBPs have no prior RNA-binding related homology or annotation (Castello et al., 2013). To expand the catalog of RBPs beyond ones with sequence or structure homology, several experimental approaches have been undertaken. In budding yeast *Saccharomyces cerevisiae*, protein microarrays and affinity purification followed by mass spectrometry identified a significant number of previously unannotated RBPs, including unconventional RNA interactions for a number of known enzymes (Klass et al., 2013; Scherrer et al., 2010; Tsvetanova et al., 2010). More recently, UV cross-linking of proteins to RNA followed by oligo(dT) capture and mass spectrometry (referred to as interactome capture) has been used to identify proteins which bind directly to mRNA. This general approach has identified nearly 900 RBPs in human HeLa cells, 729 RBPs in human HuH-7 cells, 382 nuclear RBPs in human K562 cells, 555 RBPs in mouse embryonic stem cells, 678 RBPs in yeast, and 523 in drosophila embryos (Beckmann et al., 2015; Castello et al., 2012; Conrad et al., 2016; D et al., 2016; Kwon et al., 2013; Matia-Gonzalez et al., 2015; Sysoev et al., 2016; Wessels et al., 2016). However, a major limitation of these studies is that the use of oligo(dT) beads to isolate proteins bound to polyadenylated mature messenger RNAs is biased against the identification of proteins that interact with premature, unspliced and unpolyadenylated RNA.

As many RBPs interact co-transcriptionally to influence splicing, 5' and 3' end formation, and nuclear localization, for both unpolyadenylated pre-mRNAs and non-coding RNAs, it is important to adopt approaches to extend the RBP catalog to include RBPs not identified by polyadenylated mRNA interactome capture techniques. One attractive approach is to examine protein-protein interaction (PPI) networks for both nuclear and cytoplasmic RBPs and to discriminate between RNA-mediated and direct protein-protein interactions. At a small scale, quantitative mass spectrometry-based proteomics, in combination with RNase treatments, has been successfully utilized to identify RBPs that are concurrently bound to the same RNA substrates (Flury et al., 2014; Klass et al., 2013).

Here, we propose an RNase-coupled proteomics strategy to identify proteins with RNA binding activity by prioritizing proteins that interact with multiple known (or annotated) RBPs in an RNA-dependent manner. To evaluate this concept, we first focused on RBPs localized in both the nucleus and the cytoplasm, such as the heterogeneous nuclear ribonucleoparticle proteins (HNRNPs), for which RNA binding roles are well documented and transcriptome-wide targets have been characterized (Huelga et al., 2012; Lee et al., 2015; Van Nostrand et al., 2016). We used HaloTag fusion (HT-RBP) pull-downs followed by quantitative mass spectrometry, in both untreated and RNase treated conditions, to characterize the interactomes of ten of the major HNRNPs, as well as the splicing factor (RBFOX2), and the nonsense-mediated decay (NMD) pathway regulator (UPF1). With this approach we demonstrate that interactors associating with multiple HT-RBPs are highly enriched for annotated RBPs as well candidate RBPs.

To generalize RBP discovery from PPI networks, we developed a classification algorithm Support vector machine Obtained from Neighborhood Associated RBPs (termed SONAR) to predict if a given protein of interest is a candidate RBP. Enhanced CLIP (Van Nostrand et al., 2016) performed on a subset of SONAR predictions and HT-RBP interactors

corroborated that they are bona-fide RBPs, demonstrating that SONAR successfully discovers uncharacterized candidate RBPs by leveraging protein-protein interactomes.

RESULTS

HaloTag fusion pull-down of 12 RBPs followed by quantitative proteomics identifies RNA-dependent interactors

We performed a small-scale proteomic study with and without RNase treatment to identify proteins that interact with known RBPs in an RNA-dependent manner. Specifically, we performed an unbiased quantitative proteomic study to systematically and comprehensively identify proteins in complex with 12 canonical RBPs (10 hnRNP proteins, the splicing factor RBFOX2 and nonsense-mediated decay factor UPF1) using HaloTag (HT) technology. The rapid binding kinetics of the HaloLink sepharose-based resin decreases false positives by minimizing resin exposure time to lysate, and covalent capturing allows for transient interactors to be purified by prohibiting diffusion of the primary capture target from the resin (Daniels et al., 2012; Deplus et al., 2013). Plasmids encoding a full-length open reading frame (ORF) of each RBP fused to a HaloTag (HT-RBP) were transfected in HEK293T cells in biological replicates. Confocal microscopy performed on cells expressing these HT-RBPs labeled with fluorescent TMR HaloTag ligand revealed nuclear and cytoplasmic localization (Figure S1). As these proteins are known to interact with RNA substrates in both the nucleus and cytoplasm, we expect that our fusion pull-down approach will identify interacting proteins that associate with both nascent and mature mRNAs.

Cells were lysed 24 hours post-transfection, and half of the lysates were left untreated, while the remaining half were subjected to stringent RNase digestion (Figure 1A). Briefly, protein complexes were covalently captured on HaloLink resin, subjected to washing to remove non-specific interactions and the remaining protein interactors were eluted in urea, and analyzed by silver staining (Figure S2A). Following Multidimensional Protein Identification Technology (MudPIT) mass spectrometry of eluted interactors, normalized spectral abundance factors (NSAF) values were calculated (Figure 1B), and these values were used to group replicates by hierarchical clustering (Figure S2B). Outlier experiments were removed and the remaining replicate experiments were averaged together. Within each experiment, the distribution of \log_2 enrichment scores compared to the control experiment was computed for each protein interactor (HNRNPF is presented as an example in Figure 1B, and the others in Figure S2C). The set of significantly enriched protein interactors within each experiment was defined using a Z-score. Specifically, protein interactors were significantly enriched for interacting with a given HT-RBP if they had a score greater than 1.5 standard deviations (σ) from the mean (μ) of the positive values in the distribution (red dashed line, Figure 1B). Interactor enrichment scores were computed for both the RNase-treated and untreated conditions and a total of 3853 significantly enriched interactors were detected (Table S1). For each HT-RBP, the enrichment scores for protein interactors in the untreated sample were compared to the RNase-treated sample (Figure S2D). Expectedly, the HT-RBP was often the most enriched protein precipitated in the experiments (red dots Figure S2D). Analysis of all enriched interactions for all HT-RBP baits revealed that our approach detects approximately similar numbers of RNA-independent and RNA-dependent

interacting proteins. Indeed, nearly half of all significant protein interactions with HT-RBPs are lost with RNase treatment (Figure 1C).

We found that RNase treatment depleted known RNA-dependent interactors, such as the exon-junction complex proteins (Figures 1D and 1E), which displayed enriched interactions in the untreated sample and indicates RNase-sensitive (RNA-dependent) interactions with HT-HNRNPs. In contrast, we observe known RNA-independent interactions between HT-UPF1 and the other UPF components of the EJC complex, which confirms that our experimental conditions and computational approach reliably distinguishes known RNA-independent versus RNA-dependent protein-protein interactions. Similarly, many of the interactions between HT-RBP baits and the components of the SF3B splicing complex were RNA-dependent (Figure 1D).

We performed a Gene Ontology analysis on the RNA independent interactors for each HT-RBP (Table S2). Terms associated with RNA processing are among the most statistically significantly enriched GO terms identified for HT-RBP RNA-independent interactions (Figure 1F). We also observed significantly enriched ontology terms unique to particular HNRNPs. For example, HNRNPU uniquely interacts in an RNA independent manner with E3 ubiquitin-protein ligases and proteasome subunits, HNRNPH uniquely interacts with factors involved in Golgi vesicle transport, and HNRNPF and RBFOX2 interact with unique subsets of nuclear encoded mitochondrial proteins, including respiratory electron transport chain components (Figure 1F). In summary, HaloTag fusion pull-downs of canonical RBPs successfully identified thousands of RNA-dependent and independent interactors in human cells.

RNA-dependent interactors that interact with multiple RBPs are frequently RBPs

We reasoned that the more HT-RBPs a protein interacts with, the more likely it is an RBP itself, particularly if the interactions are RNA-dependent. In fact, for both RNA-dependent and independent interactors, we observed a strong correlation between the number of HT-RBP baits that a given protein interacts with and the likelihood that the protein is itself an annotated RBP. Nearly half of all RNA-dependent and independent interactors associate with a single HT-RBP bait and a decreasing number of proteins interact with an increasing numbers of baits (Figure 2A, B). However, as the number of interacting HT-RBPs increases, the fraction of annotated RBPs represented in the set of interactors also increases, such that ~60% of interactors that interact with 5 or more HT-RBPs are annotated RBPs (Figure 2C, D). This effect is even stronger for RNA-dependent interactions. 100% of RNA-dependent interactors with more 10 or more HT-RBPs are known RBPs (Figure 2C), compared to ~70% of RNA-independent interactors (Figure 2D). Based on these analyses, we define proteins that interact with 5 or more HT-RBPs “super interactors” (SI). The ~12–20-fold enrichment of annotated RBPs within the SI proteins is notable as known RBPs currently comprise ~5–8% of all annotated protein-coding genes. Higher isoelectric points are characteristic of annotated RNA binding proteins (Castello et al., 2012), and when we compared the distribution of isoelectric points for RNA-dependent and independent SI proteins, we see average higher isoelectric points for RNA-dependent SI proteins (Figure 2E), further supporting an enrichment in RNA binding activity in this group of proteins. Our

findings support our hypothesis that RNA-dependent interactors of multiple RBPs are enriched for known or candidate RBPs.

Protein-protein interaction networks allow global prediction of RBPs

Based on our observations that the more HT-RBPs a given protein interacts with, the more likely it is an RBP, we developed a classification algorithm utilizing a linear Support-vector-machine Obtained from Neighborhood Associated RBPs (SONAR) to calculate an RBP classification score (RCS) for protein baits of interest from large-scale PPI datasets. First, we developed a list of 1787 human annotated RBPs by combining the list of 1072 by Sundararaman et al. (Sundararaman et al., 2016) and 1542 by Gerstberger et al. (Gerstberger et al., 2014) (Table S3). Next, we generated “neighborhoods” for all proteins in the human BioPlex PPI dataset, which consists of dozens of thousands of interactions (Huttlin et al., 2015). Proteins (RBP and non-RBPs) are represented as nodes in these neighborhoods with edges connecting two proteins if they interact (Figure 3A). Interestingly, known RBPs in the BioPlex dataset exhibit statistically significant differences from non-RBPs in terms of network properties. RBPs tend to have higher degree centrality (the number of edges incident upon a node; $p < 10^{-15}$) and higher closeness centrality (distance from the node to all others, $p < 10^{-11}$). RBPs and non-RBPs have approximately similar betweenness centrality (i.e. the number of times the node acts as a bridge along the shortest path between two other nodes, $p < 0.002$) compared to non-RBPs (Figure S3A). Within each neighborhood, the fractions (to mitigate differences in degree and closeness centralities) of annotated RBPs were extracted as features for the SVM classifier (Figure 3B). The classifier was trained using annotated RBPs as positively labeled examples and other (presumably non-RBP) proteins as negatively labeled examples. Performance was evaluated using 10-fold cross-validation. The areas under the receiver-operating-curve (ROC-AUC) and under the precision-recall-curve (PRC-AUC) illustrate the high classifier sensitivity (~ 0.7) and specificity (~ 0.9) of the respectively (Figures 3C and 3D). We obtained different definitions of RBPs based on interactome capture (Castello et al., 2012), Baltz (Baltz et al., 2012), Beckmann (Beckmann et al., 2015)), manual curation (Gene Ontology or GO-annotated), sequence composition alone (RBD-defined), and some combination of experimentally defined and manual curation (Gerstberger et al., 2014) (Figures 3E and S3B). Importantly, the performance of SONAR, as measured by the percentage recall, is robust to different definitions of RBPs (Figures 3E and S3B). Curiously, the RBD-defined set consisting of the largest set (2551) of proteins predicted to associate with RNA based solely on sequence composition appears to have a low fraction of SONAR-predicted RBPs, unlike the other defined lists. This strongly suggests that sequence composition alone may yield a high false positive rate when annotating RBPs. We also observed that 30% of known transcription factors are predicted by SONAR to have RNA-association activity (Figure 3E).

The RCS value for each held-out protein in the BioPlex data was determined as the mean classifier score (or output) over 10 iterations of cross-validation (Table S4). The RCS distribution for annotated RBPs (that were not used in training of SONAR) was statistically significantly higher ($p < 10^{-16}$, Kolmogorov Smirnov two-tailed test) compared to non-annotated RBPs, as well as the full set of proteins in the BioPlex data. In fact, $\sim 70\%$ of annotated RBPs had an RCS of greater than our user-defined SONAR threshold of 0.79,

whereas only 10% of all unannotated RBPs had positive scores (Figure 3F). Interestingly, 30% of the complete set of HaloTag interactors (HTI) that were also present in the BioPlex dataset scored above 0.79 (Figure 3F). The complete set of HaloTag-RBP-super interactor proteins had a score distribution similar to that for annotated RBPs. Pertinently, SI proteins that were extracted from the RNA-dependent (HT-RD-SI) interactions had significantly higher RCS values than SI proteins obtained from the RNA-independent (HT-RI-SI) interactors (Figure 3F).

Next we evaluated attributes that reflect biophysical properties thought to be associated with known RBPs (Castello et al., 2012; Lunde et al., 2007). We compared annotated RBPs and candidate proteins previously unannotated as RBPs but predicted by SONAR to be RNA-associated ($RCS > 0.79$) (blue lines, Figure S4). In comparison to all bait proteins in the BioPlex network (dashed lines, Figure S4A), we found that SONAR predicted candidate RBPs exhibited higher isoelectric points ($p < 10^{-6}$, Figure S4A), higher proportions of residues in disordered regions ($p < 10^{-5}$, Figure S4B), and higher amino acid compositions reflective of RNA binding ($p < 10^{-5}$, Figure S4C). These features were also similar to annotated RBPs (red lines, Figure S4). We did not detect any appreciable differences in low-complexity regions (Figure S4D) or total protein size (Figure S4E). In summary, SONAR leverages publically available PPI networks to predict proteins that have attributes of RBPs.

SONAR identifies hundreds of previously unannotated RBPs across multiple species

The flexibility of SONAR was next evaluated on available yeast (*Saccharomyces cerevisiae*) and fly (*Drosophila melanogaster*) protein-protein interactome datasets. Training and testing for the classifiers followed the same procedure as for the human dataset. Interestingly, we found that SONAR performed very well ($AUC > 0.8$) on yeast (Figure 4A) and fly (Figure 4B) PPI datasets generated by affinity precipitation followed by mass spectrometry (AP-MS), similar to BioPlex. Performance metrics are provided for human, yeast and fly in Figure S3C (chosen to achieve a false discovery rate of 10%). However, SONAR performs weakly using PPI datasets from yeast two-hybrid (Y2H) methodology (Figure S3D). As Y2H relies on bait and prey proteins to interact, bringing together the activation domain and binding domain of a transcription factor to activate a reporter gene, we expect that when these bait and preys are RBPs, they are not directly bound to RNA molecules. This result supports the concept that AP-MS experiments identifies endogenous RBP interactors as they are more likely to be in their normal cellular context of binding RNA. As with RCS distributions derived from the human BioPlex network, the yeast (Figure 4C) and fly (Figure 4D) RCS distributions for annotated RBPs was statistically significantly higher ($p < 10^{-16}$ Kolmogorov Smirnov two-tailed test) compared to non-annotated RBPs, as well as the full set of proteins in the BioGrid data.

SONAR candidate RBPs overlap with 31% of annotated yeast (Figure 4E) and 17% of annotated fly RBPs (Figure 4F). Here we defined annotated yeast and fly RBPs using the union of human orthologs and proteins identified from interactome capture studies (Figures S5A and S5B for yeast and S5C for fly). Performance for SONAR was also robust when trained using RBP definitions solely derived from interactome-capture or only from computationally-defined orthologs (Figure S5E). Importantly, 350 human, 290 fly and 279

yeast proteins that were predicted by SONAR as candidate RBPs had assignable orthologs, with 142 SONAR predicted candidate RBPs conserved across all three species (Figure 4G). These conserved candidate RBPs are involved in protein-folding, chromosome organization, and nuclear import, and include HT-RBP SI proteins (Table S4). In all, these results demonstrate that SONAR identifies previously unannotated and conserved RBPs across multiple species without explicit sequence and structure homology information.

Candidate RBPs are enriched for proteins with zinc finger and DNA binding domains

The set of BioPlex proteins that SONAR predicts as RBPs ($RCS > 0.79$) overlaps with ~70% of the list of annotated RBPs contained in BioPlex and ~80% of HT-RD SI proteins from the HaloTag PPI network (Figure 5A). There are 998 SONAR candidate RBPs that have no previous annotation as RBPs, 10 of which are HT-RD-SI RBP candidates (Figure 5A). We performed gene ontology (GO) analysis to determine what categories of proteins are enriched in SONAR predicted RBPs aside from the expected RNA processing categories. Intriguingly, for previously unannotated SONAR predicted RBPs, we found significant enrichment for GO terms involved in transcription, chromatin modification, chromosome organization, DNA binding, and nucleosome assembly (Figure 5B, Table S5), and protein domains contained by BioPlex candidate RBPs are significantly enriched for DNA binding domains, with zinc finger binding domains being the most enriched (Figure 5C, Table S5). We also found that SONAR predicted RBPs overlapped with 27% of the newly identified non-canonical RBPs termed “enigmRBPs” (Beckmann et al., 2015), as well as ~68% of recently reported DNA-RNA binding proteins (Conrad et al., 2016) (Figure S5D). In fact, the concordance with the latter study is also consistent with our recent data generated as part of the ENCODE consortium. We identified 10 previously annotated DNA binding proteins that SONAR predicted as RBPs ($RCS > 0.79$), which interact directly with endogenous RNA by eCLIP (Table S5, publicly available datasets at <https://www.encodeproject.org>). We believe that the large set of SONAR predicted RBP candidates (998) would significantly expand the existing repertoire of RNA-associated proteins upon further experimental confirmation.

CLIP and eCLIP validate RNA binding activity and discover RNA targets of candidate RBPs

We next sought to determine SONAR’s ability to predict RNA *in vivo* binding activity for a selection of HT-RD-SI proteins with SONAR RCS values between -1.3 and $+1.7$. We transiently transfected HEK293T cells with plasmids expressing V5 tagged fusions of a selection of candidate RBPs. Cells were UV irradiated and RBP-RNA complexes were immunoprecipitated using commercially available V5 antibody and T4 PNK radiolabeled, separated on an SDS-PAGE gel and transferred to a nitrocellulose membrane for Western blot and phosphor-image analysis. Negative controls were either non-transfected, or expressed V5-tagged GFP, while positive controls expressed the V5-tagged RBPs HNRNPM and HNRNPUL1. We found that 7 out of 8 V5-tagged RBPs tested efficiently bound radiolabeled RNA in a UV-crosslinking dependent manner (Figure S6A, Supplemental Information). In contrast, we did not detect radiolabeled RNA immunoprecipitated from cells expressing V5-tagged GFP in either condition, or from non-transfected cells. This assay validates SONAR predicted ($RCS > 0.79$) at 100%, but we note here that 3 HT-RD-SI proteins (AIFM1, VIM and SMU1) with low SONAR scores also detectably bind RNA

above negative controls by this assay, indicating a false negative rate for SONAR (Figure S6A).

Next we determined the transcriptome-wide binding sites of a subset of these RBP candidates, and included the zinc-finger protein of unknown function ZNF184, using enhanced CLIP. Briefly, HEK293T were subjected to UV-mediated crosslinking, lysis and treatment with limiting amount of RNase, followed by immunoprecipitation (IP) of protein-RNA complexes using commercially available antibodies that interrogate the endogenous proteins. RNA fragments protected from RNase digestion by RBP occupancy were subjected to 3' RNA linker ligation, reverse-transcription and 3' DNA linker ligation to generate eCLIP libraries for high-throughput Illumina sequencing. The improved efficiency of eCLIP enabled the generation of a Size-Matched Input (SMInput) library for each biological sample, in which 2% of the pre-immunoprecipitation sample was subjected to identical library generation steps including ribonuclear protein complex size-selection on nitrocellulose membranes. In total, 42 eCLIP (including SMInput) libraries were sequenced to ~7 million reads each, of which ~30–70% mapped uniquely to the human genome (Table S6).

We demonstrate that these candidate RBPs bind RNAs in different genic regions (Figure 6A) and recognize different binding motifs (Figure 6B). The RAN-GTPase activating protein RANGAP1 showed preferential enrichment in 5' UTRs and intron-less genes (Figures 6A, 6C). RANGAP1 localizes at the nuclear periphery and is important for import of cargo through nuclear pore complexes. Disrupted nuclear import caused by RANGAP1 inhibition through interaction with C9orf72 hexanucleotide repeat expansion RNA (HREs) is involved in ALS pathology (Zhang et al., 2015). Interestingly, we find that the motif (CGGCGG) enriched by RANGAP1 eCLIP is similar to the G₄C₂ pattern responsible for HRE G-quadruplexes (Figure 6B). The mitotic spindle organizing protein NUMA1 preferentially bound to intronic regions including its own pre-mRNA (Figure 6A, D). The ring-type zinc finger protein RNF219, previously implicated in Alzheimer's disease, preferentially bound 3'UTRs (Figure 6A, E). The zinc finger containing DNA binding protein ZNF184 had enriched clusters over all types of transcript regions, including introns such as the distal intron of the CENPM gene (Figures 6A, F). In all, SONAR-predicted RBPs indeed interact with hundreds to thousands of enriched binding sites in the human transcriptome.

CLIP cluster discovery was performed using the CLIPper (Lovci et al., 2013) algorithm. SI proteins that had low RCS values, such as AIFM and VIM, had fewer significantly enriched clusters (Figure S6C). Predicted RBPs with high RCS values such as NUMA1, RANGAP1, RNF219, and ZNF184 had far more enriched clusters and bound transcripts in specific regions (Figures 6A, S6C). Interestingly, SONAR RCS values appear positively correlated with the fraction of interactions with HT-RBPs that are RNA-dependent ($R^2=0.72$) (Figure S6B), and with the number of identified eCLIP clusters ($R^2=0.42$) (Figure S6C). Overall these findings demonstrate that the SONAR algorithm is effective at predicting bona-fide RBPs.

DISCUSSION

The increasing importance of RNA binding proteins in development and disease has accelerated the need to comprehensively, rapidly and accurately identify new components of ribonuclear particles in a variety of organisms. Previous computational approaches to define the complete catalog of RNA binding proteins rely heavily on the presence of RNA binding domains (Baltz et al., 2012; Gerstberger et al., 2014). These approaches lack sensitivity given the growing number of interactome capture studies that identify RBPs that lack characterized RBDs (Beckmann et al., 2015; Conrad et al., 2016). However, interactome capture studies of polyA-selected mRNPs may not identify proteins interacting with nascent or pre-mRNA and interactome capture followed by mass spectrometry is a technically demanding technique, precluding its adoption for the vast majority of researchers.

In this study, we have developed a simple yet powerful computational approach to discover RNA binding proteins based on the realization that without RNAase treatment, pull-down of RNA-associated proteins enriches for interactors that are directly bound to RNA. Therefore, the higher the number of interacting RBPs, the more likely a previously unannotated protein is itself an RBP. As a pilot study, we have independently identified RBP-interactomes for a diverse set of RBPs, namely a splicing factor (RBFOX2), a decay regulator (UPF1) and 10 major HNRNP proteins using a HaloTag purification approach followed by mass spectrometry analysis. Together, these 12 proteins are localized in the nucleus and cytoplasm, and are ideal for identifying known and previously unannotated RBPs bound to both nuclear and cytoplasmic RNA substrates. By comparing RNase treated and untreated conditions, we were able to uncover thousands of RNA-dependent and independent interactions. Interactors include many unexpected proteins involved in diverse biological processes including oxidative phosphorylation, proteasome function, energy production, mitochondrial organization, and Golgi vesicle transport. Importantly, this resource confirms that proteins that interact with multiple HT-RBPs are enriched for RBPs.

As there are many large-scale PPI datasets available that are not subject to RNase treatment, we developed an algorithm termed SONAR that leverages the neighborhood of RBP interactors in PPI networks to identify unannotated RBPs. SONAR does not rely on sequence or structural homology or modeling which historically prevents the identification of proteins that have RNA binding activity through yet unknown mechanisms. Instead, SONAR relies on experimentally determined interaction information from thousands of independently generated affinity purification experiments which in aggregate is more stringent than replicates of a few to a dozen precipitation experiments. Furthermore, as we demonstrated, SONAR can be readily applied to PPI datasets from multiple organisms, which led to our identification of previously unannotated, evolutionarily conserved RBPs. Interestingly, many of these candidate RBPs have zinc finger binding domains and are involved in processes such as transcription, chromosome organization and chromatin modification, in support of recent predictions that there are likely RNA binding roles for DNA binding proteins (Conrad et al., 2016). While it is possible that this observation is the result of direct protein-protein interactions between DNA-binding protein baits and nuclear RBPs, our results also suggest the possibility that many DNA binding proteins likely also interact with RNA during transcription or RNA processing. The second possibility is

consistent with our own eCLIP studies performed as part of the ENCODE consortium, which show that 10 high SONAR-scoring proteins previously reported to have DNA-binding and transcription factor activity can also interact directly with endogenous RNA targets.

Lastly, we validated SONAR predictions of RNA binding activity by eCLIP-seq for a set of RBP candidates. Surprisingly, for 6 candidate eCLIP validations, higher SONAR scores correlated with higher numbers of eCLIP-identified binding sites, with VIM and AIFM libraries having the fewest peaks above size-matched input controls, and NUMA1, RANGAP1 and ZNF184 having many more clusters enriched above size-matched input controls. In retrospect, this interesting observation is consistent with our features utilized by SONAR. It is reasonable to imagine that the higher the fraction of known RBPs within the neighborhoods of a candidate RBP, the higher the number of different RNA molecules bound by various RBPs (hence more binding sites) that the protein should interact with.

Among our validated examples of previously unannotated RNA-associated proteins, the ring-type zinc finger protein RNF219 preferentially bound 3'UTRs, which is intriguing in light of its reported protein-protein interactions with components of the deadenylation machinery (Hein et al., 2015; Li et al., 2014), and it will be interesting in the future to determine whether RNF219 affects stability of its RNA targets. The mitotic spindle organizing protein NUMA1 preferentially bound introns, and considering the strong RNA independent interaction of NUMA1 with UPF1 and RBFOX2 baits, it is possible that NUMA1 plays a role in alternative splicing or nonsense mediated decay of its targets. The RAN-GTPase activating protein RANGAP1 preferential bound many 5' UTRs. Disruption of RANGAP1 mediated mRNA export and localization could be a possible mechanism underlying C9orf72 HRE related ALS pathology (Zhang et al., 2015), and in future experiments, it is important to determine if C9orf72 HRE disruption of RANGAP1 RNA binding may lead to mRNA export/import inhibition and altered stability or translation for RANGAP1 mRNA targets. The DNA binding protein ZNF184 bound transcripts in introns, UTRs and coding sequences. Given the recent finding that mRNA fate can be determined by promoter regions alone (Zid and O'Shea, 2014), one interesting explanation of this phenomenon may be that RNA binding activity of promoter specific transcription factors and chromatin remodelers influences transcript localization, and the group of transcription factors predicted to bind RNA in this study could serve as valuable candidates for testing this hypothesis.

In conclusion, the SONAR algorithm is a powerfully simple and effective method for the *de novo* identification of candidate RBPs. As with many methods, the strength of SONAR's approach, which lies on its ability to leverage large-scale PPI datasets, is also a potential weakness. Firstly, SONAR relies on the coverage of these PPI networks, meaning not all proteins have been subjected to affinity precipitation followed by mass spectrometry. Secondly, SONAR relies on the quality of the independent datasets within these PPI networks. Finally, SONAR requires definitions of known or annotated RBPs to generate its scores for uncharacterized proteins. Fortunately, as more interactome capture datasets are generated and more SONAR-predicted RBPs are validated, we can build an increasingly rigorous dataset of known RBPs. Also, efforts to improve the quality and coverage of PPI datasets are underway. Nevertheless, in this study, the RBP interactomes presented here

represent a valuable first step in building a comprehensive resource for the identification of RNA-dependent and independent interactors of RBPs. Also, future studies verifying the large set of classifier predicted RBPs would expand the list of RBPs to include many unexpected classes of proteins involved in diverse biological processes.

EXPERIMENTAL PROCEDURES

HaloTag Affinity Purification

HaloTag affinity purification was performed as follows, with and without RNase treatment in replicate (2 μ l RNase A; A797C 4mg/ml for each 12 million cell pellet lysate). Protein complexes were covalently captured on HaloLink resin pre-equilibrated in resin wash buffer (TBS and 0.05% IGEPAL CA-640; Sigma) for 15 min at 22°C with rotation. Resin was then washed five times with wash buffer to remove non-specific interactions, and protein interactors were eluted in 8M Urea. Additional control experiments with HaloTag only were also performed with and without RNase digestion in replicate. Affinity purified complexes were then analyzed by Multidimensional Protein Identification Technology (MudPIT) as previously described (Florens and Washburn, 2006).

Mass Spectrometry data analysis

Spectral counting was performed and normalized spectral abundance factors determined (dNSAF values) as previously described (Zhang et al., 2010). False positives/contaminants were removed by comparing to control experiments. Any potential protein interactor with a dNSAF value that was not 2 \times greater than its value in any of the control corresponding experiments was assigned a 0 value.

For equivalent comparison to the RBP RNA targets, the lists of protein-protein interactors were filtered for those targets that are expressed at a detectable RNA level in HEK293T cell RNA-seq data (Huelga et al., 2012). Replicates were clustered, outlier datasets (UPF1_RNASE_2, M_RNASE_1, H_RNASE_1, D0_1) were removed, and the remaining replicates were averaged together. An enrichment score was calculated for each protein interactor by calculating a log₂ ratio of the dNSAF value in the RBP purification experiment compared to the Control experiment:

$$\text{Enrichment Score} = \log_2[(\text{dNSAF RBP} + 0.000001)/(\text{dNSAF Control} + 0.000001)]$$

A pseudocount of 0.000001 representing the lowest measurable dNSAF value was added within the enrichment score to prevent losing targets with zero values in the Control experiment (divide by zero errors). The final set of protein-protein interactors was determined using all protein interactors with enrichment greater than the $\mu - 1.5\sigma$ in each dataset.

RNA-dependent interactions were defined as: $((\text{no_RNase enrichment score}) - (\text{RNase enrichment score})) / ((\text{no_RNase enrichment score}) + (\text{RNase enrichment score})) > 0.3$. Post lysis interactions were filtered and defined as: $((\text{no_RNase enrichment score}) - (\text{RNase enrichment score})) / ((\text{no_RNase enrichment score}) + (\text{RNase enrichment score})) < -0.2$. All other interactions were defined as RNA independent.

Data Accessibility: The complete MudPIT datasets can be obtained from <ftp://massive.ucsd.edu/> using the following MassIVE accession numbers for the MOCK and hnRNPs datasets with password (KBGWY39029): MSV000079668 and MSV000079669. The corresponding ProteomeXchange accessions are: PXD003999 and PXD004000.

PPI network construction and neighborhood features generation

All the protein-protein interaction networks were represented as an undirected and unweighted graph $G = (V, E)$. Each vertex $v \in V$ refers to a protein and each edge $(v, u) \in E$ displays an interaction between two proteins, v and u . The definition of level- k is described in (Sela et al., 2012): briefly, u is v 's level- k neighbor when there exists a path of length k from v to u . It should be noted that there may be more than one path between v and u , thus u can be v 's level- k_1 and level- k_2 ($k_1 \neq k_2$) simultaneously. Additionally, if there are n ($n > 1$) paths of k from v to u , the protein u will be counted k times in the feature calculation of level- k neighbors. In this study, we only consider the given protein's first three levels of neighbors (shown in Figure 3A), and an example of the feature generation process is displayed in Figure 3B. With the three features (calculations described in Supplemental information), we obtain a feature vector for each protein of the PPI network. We constructed a human protein-protein interaction (PPI) network using interaction data from the BioPlex project (Huttlin et al., 2015). For each protein in this network, a PPI feature was generated with the approach above. PPI networks for 2 other model organisms (*Saccharomyces cerevisiae* and *Drosophila melanogaster*) were constructed with the interaction data downloaded from BioGRID 3.4 database (Stark et al., 2006). For each species, two PPI networks were assembled with the data from affinity purification – mass-spectrometry (AP-MS) and yeast-two-hybrid experiments respectively. The 3-level PPI features were generated using the same calculation above for each protein in the PPI networks.

Classifier building and evaluation

Classifiers were trained to recognize RBPs using the SVM classification framework with a Radial Basis Function (RBF) kernel as implemented within the scikit-learn package in Python (Pedregosa et al., 2011). We divided the proteins in each PPI network into two categories, RBPs and non-RBPs, based on the RBP annotation lists as described in the main text. Each protein is represented by a feature vector and a class label (i.e. RBP or non-RBP). We trained the classifiers and evaluated their performance using Receiver Operating Curve (ROC)-Area Under the Curve (AUC) and Precision-Recall (PRC)-AUC analysis with 10-fold cross-validation. The closer the classifier's ROC curve to the upper left corner, the better the classifier's performance in both sensitivity and specificity. To train our classifiers with a balanced dataset, we oversampled RBP samples in the training step. During testing, each unseen protein obtained a classification score from a trained classifier. We generated classification scores for every protein in the PPI network during the cross-validation training/testing process and repeated 10 times to get an average score for each protein.

eCLIP

eCLIP was performed as previously described (Van Nostrand et al., 2016). Briefly, for each experiment, 20 million HEK293T cells were UV crosslinked (254 nm, 400 mJ/cm²) and lysed on ice, and lysates were sheared with RNase I (Ambion). RBP-RNA complexes were

immunoprecipitated with antibodies (see Supplemental Information) specific to candidates of interest (Protein G sheep anti-rabbit Dynabeads), and immunoprecipitated material was stringently washed. Dephosphorylation with FastAP (Thermo Fisher) and T4 PNK (NEB), are followed by on-bead ligation of barcoded RNA adapters to the 3' end (T4 RNA Ligase, NEB). RNA-protein complexes are run on standard protein gels and transferred to nitrocellulose membranes, and the region 75 kDa (~150 nt of RNA) above the protein molecular weight is excised and proteinase K (NEB) treated to isolate RNA. Reverse transcription is carried out using Affinityscript (Agilent), followed by treatment with ExoSAP-IT (Affymetrix) to remove excess oligonucleotides. A DNA randomer adapter is then ligated to the cDNA fragment 3' end (T4 RNA Ligase, NEB). After cleanup (Dynabeads MyOne Silane, ThermoFisher), samples were first subjected to qPCR to determine final PCR cycle number, and then PCR amplified (Q5, NEB) and size selected via agarose gel electrophoresis. Samples were sequenced on the Illumina HiSeq 2500 platform as two Paired End 50bp (for N5) or 55bp (for N10) reads.

Supplementary Material

Refer to Web version on PubMed Central for supplementary material.

Acknowledgments

The authors would like to thank members of the Yeo lab, especially Olga Botvinnik and Katannya Kapeli for critical reading of the manuscript. This work was supported by grants from the National Institute of Health (HG004659, U54HG007005 and NS075449 to G.W.Y.). C.A.S.B., J.M.G., L.F. and M.P.W. are supported by the Stowers Institute for Medical Research. E.L.V.N. is a Merck Fellow of the Damon Runyon Cancer Research Foundation (DRG-2172-13). G.W.Y. is an Alfred P. Sloan Research Fellow. G.A.P. is a National Science Foundation Graduate Fellow and was partially supported by the University of California, San Diego, Genetics Training Program through an institutional training grant from the National Institute of General Medical Sciences, T32 GM008666. K.B. is a University of California President's Postdoctoral Fellow and was partially supported by the California Institute for Regenerative Medicine Training Program (CIRM).

REFERENCES

- Baltz AG, Munschauer M, Schwanhaussner B, Vasile A, Murakawa Y, Schueler M, Youngs N, Penfold-Brown D, Drew K, Milek M, et al. The mRNA-bound proteome and its global occupancy profile on protein-coding transcripts. *Molecular cell*. 2012; 46:674–690. [PubMed: 22681889]
- Beckmann BM, Horos R, Fischer B, Castello A, Eichelbaum K, Alleaume AM, Schwarzl T, Curk T, Foehr S, Huber W, et al. The RNA-binding proteomes from yeast to man harbour conserved enigmRBPs. *Nat Commun*. 2015; 6:10127. [PubMed: 26632259]
- Castello A, Fischer B, Eichelbaum K, Horos R, Beckmann BM, Strein C, Davey NE, Humphreys DT, Preiss T, Steinmetz LM, et al. Insights into RNA biology from an atlas of mammalian mRNA-binding proteins. *Cell*. 2012; 149:1393–1406. [PubMed: 22658674]
- Castello A, Fischer B, Hentze MW, Preiss T. RNA-binding proteins in Mendelian disease. *Trends in genetics : TIG*. 2013; 29:318–327. [PubMed: 23415593]
- Conrad T, Albrecht AS, de Melo Costa VR, Sauer S, Meierhofer D, Orom UA. Serial interactome capture of the human cell nucleus. *Nat Commun*. 2016; 7:11212. [PubMed: 27040163]
- D GH, Kelley DR, Tenen D, Bernstein B, Rinn JL. Widespread RNA binding by chromatin-associated proteins. *Genome Biol*. 2016; 17:28. [PubMed: 26883116]
- Daniels DL, Mendez J, Mosley AL, Ramisetty SR, Murphy N, Benink H, Wood KV, Urh M, Washburn MP. Examining the complexity of human RNA polymerase complexes using HaloTag technology coupled to label free quantitative proteomics. *Journal of proteome research*. 2012; 11:564–575. [PubMed: 22149079]

- Deplus R, Delatte B, Schwinn MK, Defrance M, Mendez J, Murphy N, Dawson MA, Volkmar M, Putmans P, Calonne E, et al. TET2 and TET3 regulate GlcNAcylation and H3K4 methylation through OGT and SET1/COMPASS. *The EMBO journal*. 2013; 32:645–655. [PubMed: 23353889]
- Florens L, Washburn MP. Proteomic analysis by multidimensional protein identification technology. *Methods Mol Biol*. 2006; 328:159–175. [PubMed: 16785648]
- Flury V, Restuccia U, Bachi A, Muhlemann O. Characterization of Phosphorylation- and RNA-Dependent UPF1 Interactors by Quantitative Proteomics. *Journal of proteome research*. 2014
- Gerstberger S, Hafner M, Tuschl T. A census of human RNA-binding proteins. *Nat Rev Genet*. 2014; 15:829–845. [PubMed: 25365966]
- Hein MY, Hubner NC, Poser I, Cox J, Nagaraj N, Toyoda Y, Gak IA, Weisswange I, Mansfeld J, Buchholz F, et al. A human interactome in three quantitative dimensions organized by stoichiometries and abundances. *Cell*. 2015; 163:712–723. [PubMed: 26496610]
- Huelga SC, Vu AQ, Arnold JD, Liang TY, Liu PP, Yan BY, Donohue JP, Shiue L, Hoon S, Brenner S, et al. Integrative genome-wide analysis reveals cooperative regulation of alternative splicing by hnRNP proteins. *Cell Rep*. 2012; 1:167–178. [PubMed: 22574288]
- Huttlin EL, Ting L, Bruckner RJ, Gebreab F, Gygi MP, Szpyt J, Tam S, Zarraga G, Colby G, Baltier K, et al. The BioPlex Network: A Systematic Exploration of the Human Interactome. *Cell*. 2015; 162:425–440. [PubMed: 26186194]
- Klass DM, Scheibe M, Butter F, Hogan GJ, Mann M, Brown PO. Quantitative proteomic analysis reveals concurrent RNA-protein interactions and identifies new RNA-binding proteins in *Saccharomyces cerevisiae*. *Genome Res*. 2013; 23:1028–1038. [PubMed: 23636942]
- Kwon SC, Yi H, Eichelbaum K, Fohr S, Fischer B, You KT, Castello A, Krijgsvelde J, Hentze MW, Kim VN. The RNA-binding protein repertoire of embryonic stem cells. *Nature structural & molecular biology*. 2013; 20:1122–1130.
- Lee SR, Pratt GA, Martinez FJ, Yeo GW, Lykke-Andersen J. Target Discrimination in Nonsense-Mediated mRNA Decay Requires Upf1 ATPase Activity. *Molecular cell*. 2015; 59:413–425. [PubMed: 26253027]
- Li S, Wang L, Fu B, Berman MA, Diallo A, Dorf ME. TRIM65 regulates microRNA activity by ubiquitination of TNRC6. *Proc Natl Acad Sci U S A*. 2014; 111:6970–6975. [PubMed: 24778252]
- Lovci MT, Ghanem D, Marr H, Arnold J, Gee S, Parra M, Liang TY, Stark TJ, Gehman LT, Hoon S, et al. Rbfox proteins regulate alternative mRNA splicing through evolutionarily conserved RNA bridges. *Nature structural & molecular biology*. 2013; 20:1434–1442.
- Lukong KE, Chang KW, Khandjian EW, Richard S. RNA-binding proteins in human genetic disease. *Trends in genetics : TIG*. 2008; 24:416–425. [PubMed: 18597886]
- Lunde BM, Moore C, Varani G. RNA-binding proteins: modular design for efficient function. *Nature reviews Molecular cell biology*. 2007; 8:479–490. [PubMed: 17473849]
- Matia-Gonzalez AM, Laing EE, Gerber AP. Conserved mRNA-binding proteomes in eukaryotic organisms. *Nature structural & molecular biology*. 2015; 22:1027–1033.
- Pedregosa F, Varoquaux G, Gramfort A, Michel V, Thirion B, Grisel O, Blondel M, Prettenhofer P, Weiss R, Dubourg V, et al. Scikit-learn: Machine Learning in Python. *J Mach Learn Res*. 2011; 12:2825–2830.
- Scherrer T, Mittal N, Janga SC, Gerber AP. A screen for RNA-binding proteins in yeast indicates dual functions for many enzymes. *PloS one*. 2010; 5:e15499. [PubMed: 21124907]
- Sela D, Chen L, Martin-Brown S, Washburn MP, Florens L, Conaway JW, Conaway RC. Endoplasmic reticulum stress-responsive transcription factor ATF6alpha directs recruitment of the Mediator of RNA polymerase II transcription and multiple histone acetyltransferase complexes. *The Journal of biological chemistry*. 2012; 287:23035–23045. [PubMed: 22577136]
- Si J, Cui J, Cheng J, Wu R. Computational Prediction of RNA-Binding Proteins and Binding Sites. *Int J Mol Sci*. 2015; 16:26303–26317. [PubMed: 26540053]
- Stark C, Breitkreutz BJ, Reguly T, Boucher L, Breitkreutz A, Tyers M. BioGRID: a general repository for interaction datasets. *Nucleic Acids Res*. 2006; 34:D535–D539. [PubMed: 16381927]
- Sundaraman B, Zhan L, Blue SM, Stanton R, Elkins K, Olson S, Wei X, Van Nostrand EL, Pratt GA, Huelga SC, et al. Resources for the Comprehensive Discovery of Functional RNA Elements. *Molecular cell*. 2016; 61:903–913. [PubMed: 26990993]

- Sysoev VO, Fischer B, Frese CK, Gupta I, Krijgsveld J, Hentze MW, Castello A, Ephrussi A. Global changes of the RNA-bound proteome during the maternal-to-zygotic transition in *Drosophila*. *Nat Commun*. 2016; 7:12128. [PubMed: 27378189]
- Tsvetanova NG, Klass DM, Salzman J, Brown PO. Proteome-wide search reveals unexpected RNA-binding proteins in *Saccharomyces cerevisiae*. *PloS one*. 2010; 5
- Van Nostrand EL, Pratt GA, Shishkin AA, Gelboin-Burkhart C, Fang MY, Sundararaman B, Blue SM, Nguyen TB, Surka C, Elkins K, et al. Robust transcriptome-wide discovery of RNA-binding protein binding sites with enhanced CLIP (eCLIP). *Nat Methods*. 2016
- Wessels HH, Imami K, Baltz AG, Kolinski M, Beldovskaya A, Selbach M, Small S, Ohler U, Landthaler M. The mRNA-bound proteome of the early fly embryo. *Genome Res*. 2016
- Zhang K, Donnelly CJ, Haeusler AR, Grima JC, Machamer JB, Steinwald P, Daley EL, Miller SJ, Cunningham KM, Vidensky S, et al. The C9orf72 repeat expansion disrupts nucleocytoplasmic transport. *Nature*. 2015; 525:56–61. [PubMed: 26308891]
- Zhang Y, Wen Z, Washburn MP, Florens L. Refinements to label free proteome quantitation: how to deal with peptides shared by multiple proteins. *Anal Chem*. 2010; 82:2272–2281. [PubMed: 20166708]
- Zid BM, O'Shea EK. Promoter sequences direct cytoplasmic localization and translation of mRNAs during starvation in yeast. *Nature*. 2014; 514:117–121. [PubMed: 25119046]

Highlights

- 1.** Quantitative proteomics shows that RNA-dependent interactors are frequently RBPs.
- 2.** SONAR identifies thousands of unannotated candidate RBPs in yeast, fly and human.
- 3.** SONAR candidate RBPs are enriched for proteins with zinc finger domains.
- 4.** Enhanced CLIP identifies transcriptome-wide targets of SONAR candidate RBPs.

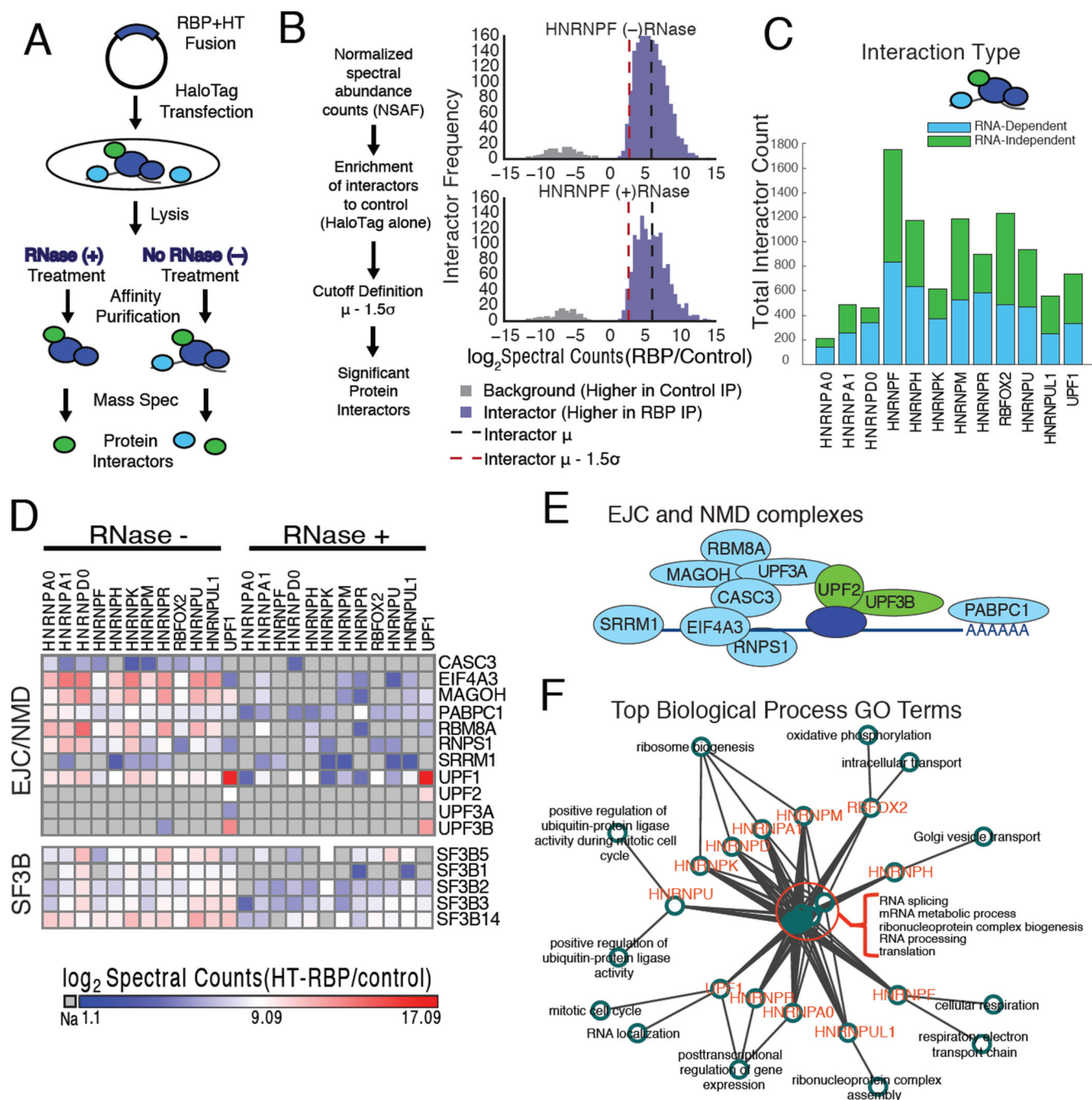


Figure 1. Identification of enriched RNA binding protein (RBP) protein interactors for HT-RBPs

A. HaloTag fusion pulldown and mass spectrometry (MS) experimental procedure. RBP-Halo Tag fusion protein constructs are transfected into HEK293T cells in replicate, cells are lysed, and half the lysate is treated with RNase. Affinity purified products are subjected to LC/MS/MS to identify protein interactors.

B. Analysis flow-chart for post processing of MS data. Normalized spectral abundance counts (NSAF) enrichment score distribution compared to control (HaloTag alone) for hnRNPF pulldown. Grey data points (enrichment <0) represent background. For enrichment

score higher in hnRNPF experiment (blue) a mean (μ , black dashed line) and standard deviation (σ) is computed. Significant interactions have enrichment score greater than 1.5 times the standard deviation (1.5σ , red dashed line).

C. Number of enriched RNA-dependent and RNA-independent interactions for all HT-RBP baits.

D. Heatmap of specific interactions displaying \log_2 fold enrichment over control for all HT-RBP baits. Interactions are grouped into NMD/EJC complexes and SF3B complexes.

E. Exon junction (EJC) and Nonsense Mediated Decay (NMD) factors. Green indicates RNA-independent interactors, and light blue indicates RNA-dependent interactors for HT-UPF1 (dark blue).

F. Gene ontology characterization of RNA-independent HT-RBP interactors. Shared and unique gene ontology terms displayed as interaction network, where red text nodes are HT-RBP baits, and central highlighted nodes are enriched terms shared by all baits (see also Table S2).

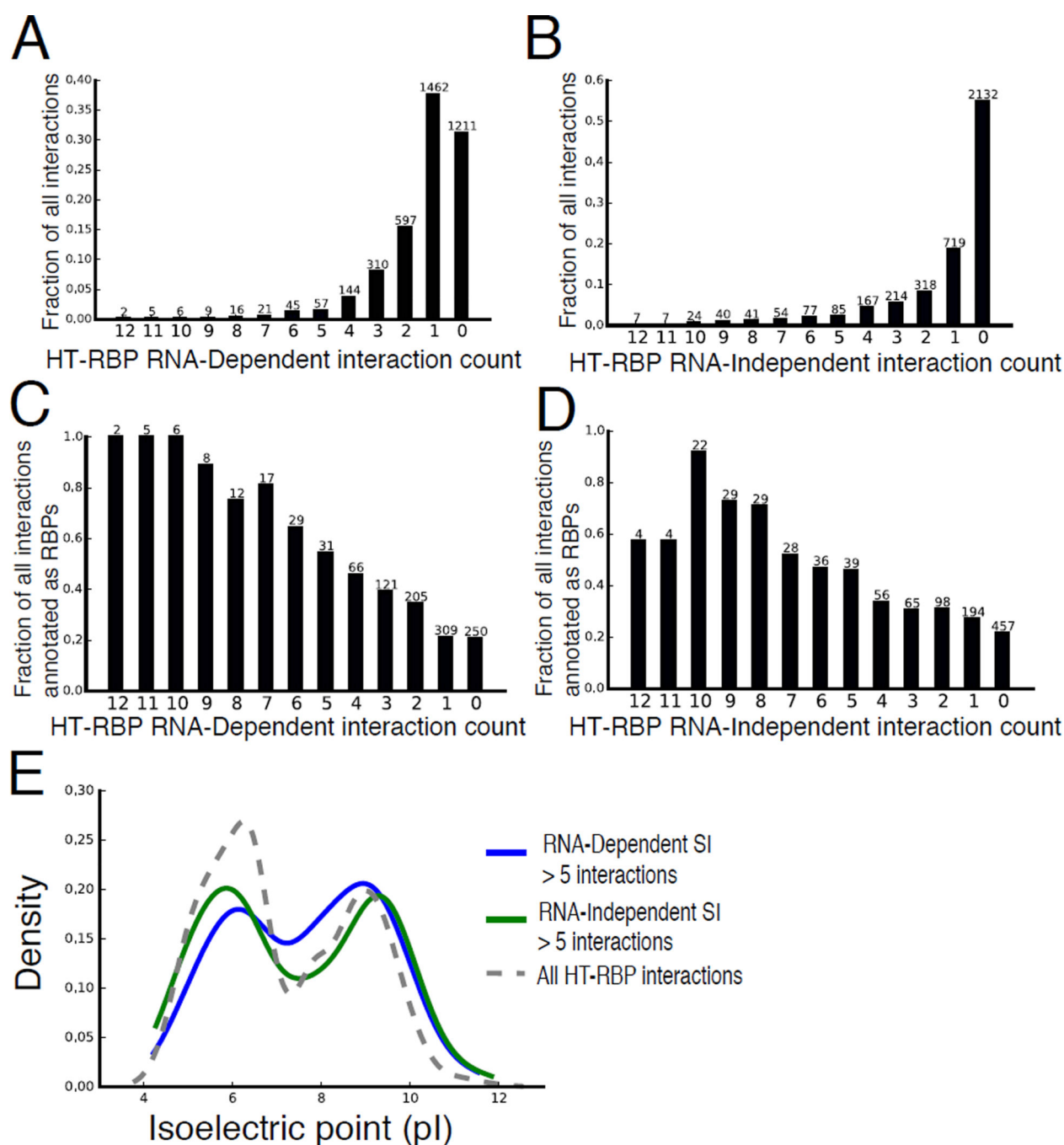


Figure 2. Super interactors are enriched for RBPs and candidate RBPs

A. Bar chart displays the fraction of all RNA-dependent interacting proteins that come up in 1 (unique interactor), and 2 to 12 HaloTag-RBP experiments (shared interactor). The number of interactors is given at the top of each bar.

B. Bar chart displays the fraction of all RNA independent interacting proteins that come up in 1 (unique interactor), and 2 to 12 HaloTag-RBP experiments (shared interactor). The number of interactors is given at the top of each bar.

C. Bar chart displays the fraction of unique (1 HT-RBP) and shared (2–12 HT-RBPs) RNA-dependent interacting proteins that are RBPs. The number of RBP interactors is given at the top of each bar.

D. Bar chart displays the fraction of unique (1 HT-RBP) and shared (2–12 HT-RBPs) RNA independent interacting proteins that are RBPs. The number of RBP interactors is given at the top of each bar.

E. Density of the calculated isoelectric points (pI) of RNA-dependent super interacting proteins (blue line), and RNA independent super interacting proteins (green line) compared to all proteins in the HT-RBP interaction set (gray dashed-line).

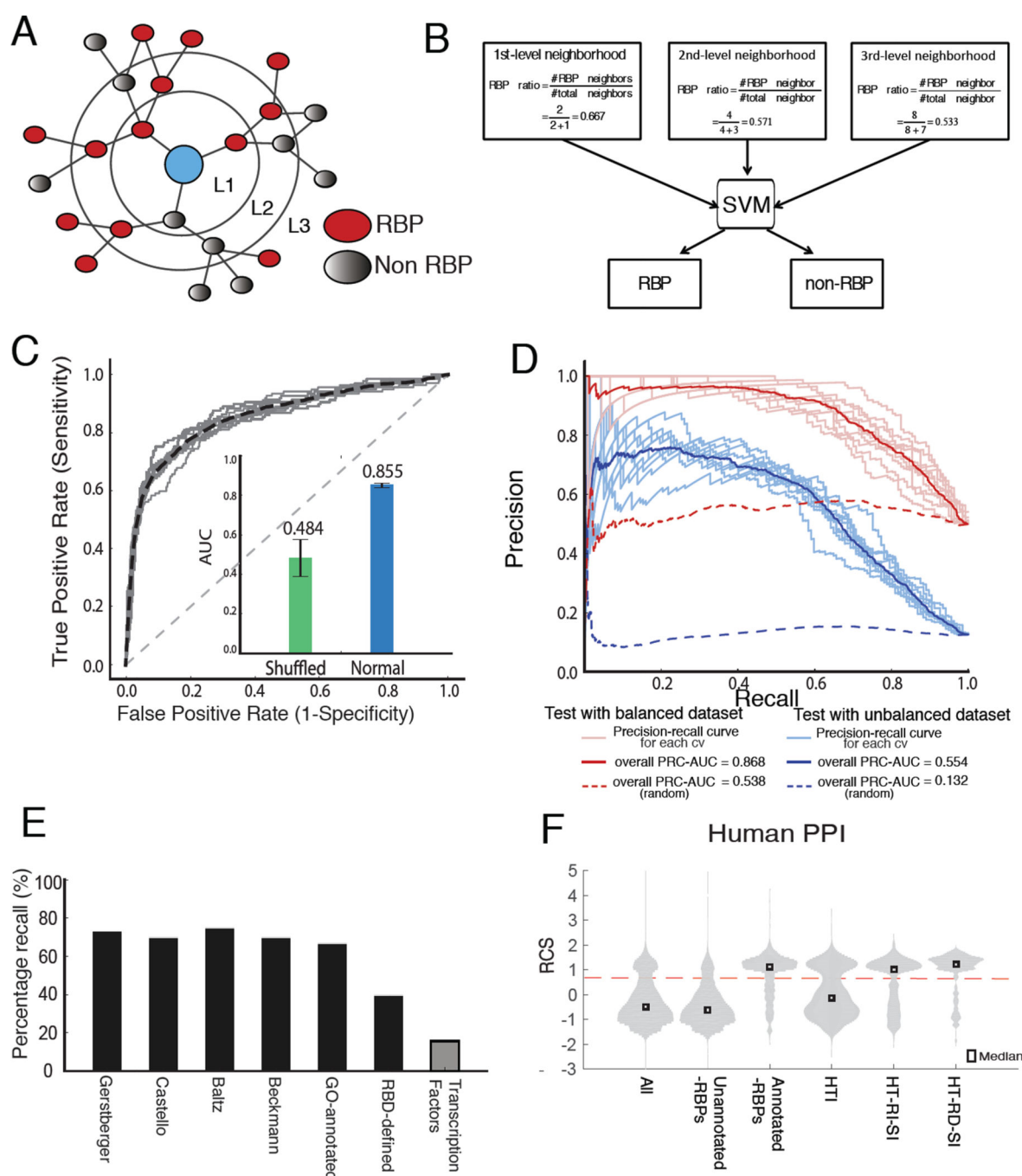


Figure 3. SONAR RBP classification approach

A. Diagram of neighborhood classification strategy. Protein of interest (POI) from a given interactome data set with its depicted neighborhood with interactions at different levels. 1st level interactions are direct interactions with POI, 2nd level interactions are interactions with 1st level neighbors, and 3rd level interactions are interactions with 2nd level neighbors.

B. Determination of RBP classification score (RCS) as described in the Methods section.

C. ROC-AUC analysis of classifier performance for human proteins from BioPlex network. Data are represented as mean \pm standard error of the mean (SEM).

D. PRC-AUC analysis of classifier performance for human proteins from BioPlex network.

E. Percent recall for SONAR trained on BioPlex PPI network for 6 RBP lists depicting different RBP annotations, and percent of annotated transcription factors (TF) predicted as candidate RBPs.

F. Violin plots of RBP classification score (RCS) distributions for all human BioPlex interactors, non-RBP interactors within BioPlex, annotated RBP interactors within BioPlex, all HT interactors, HT-RNA-independent super interactors (HT-RI-SI) and HT-RNA-dependent super interactors (HT-RD-SI) within Bioplex. The median of the distributions are denoted with a square box (see also Table S4).

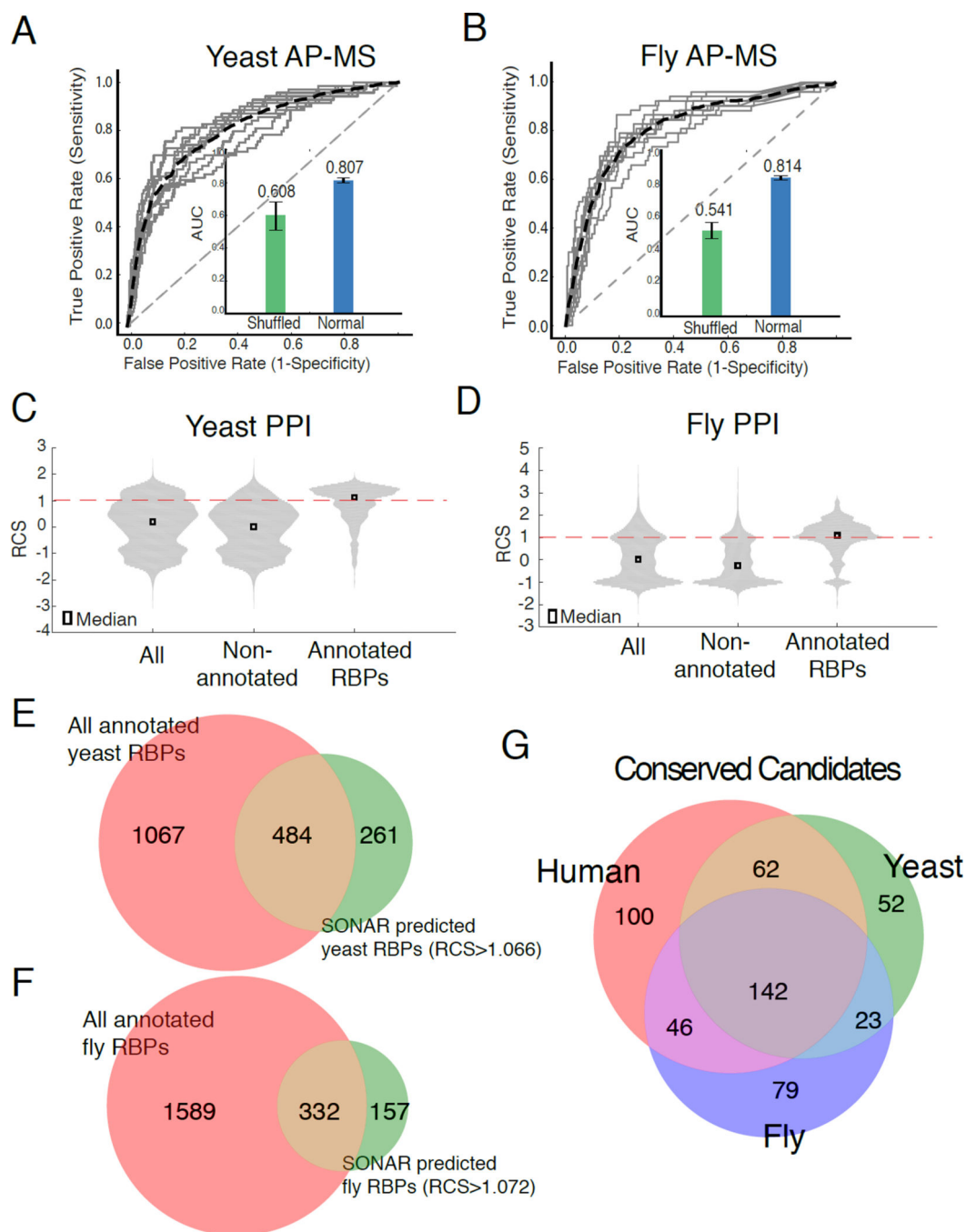


Figure 4. SONAR RBP classification scores (RCS) predict thousands of RBPs using PPI networks from multiple species

A. ROC-AUC analysis of classifier performance for yeast (*Saccharomyces cerevisiae*) proteins from BioGrid network. Data are represented as mean \pm standard error of the mean (SEM).

B. ROC-AUC analysis of classifier performance for fly (*Drosophila melanogaster*) proteins from BioGrid network. Data are represented as mean \pm standard error of the mean (SEM).

C. Violin plots of RBP classification score (RCS) distributions for all yeast BioGrid interactors, non-RBP interactors and annotated RBP interactors within yeast BioGrid interactors.

D. Violin plots of RBP classification score (RCS) distributions for all fly BioGrid interactors, non-RBP interactors and annotated RBP interactors within BioGrid interactors.

E. Venn diagram showing overlap between all annotated yeast RBPs and SONAR predicted yeast RBPs (RCS>1.066, threshold for false positive rate 0.1).

F. Venn diagram showing overlap between all annotated fly RBPs and SONAR predicted fly RBPs (at RCS>1.072, threshold for false positive rate 0.1).

G. Venn diagram showing overlap between conserved high RCS scoring predicted RBPs in human (light red), in yeast (green), and in fly (purple).

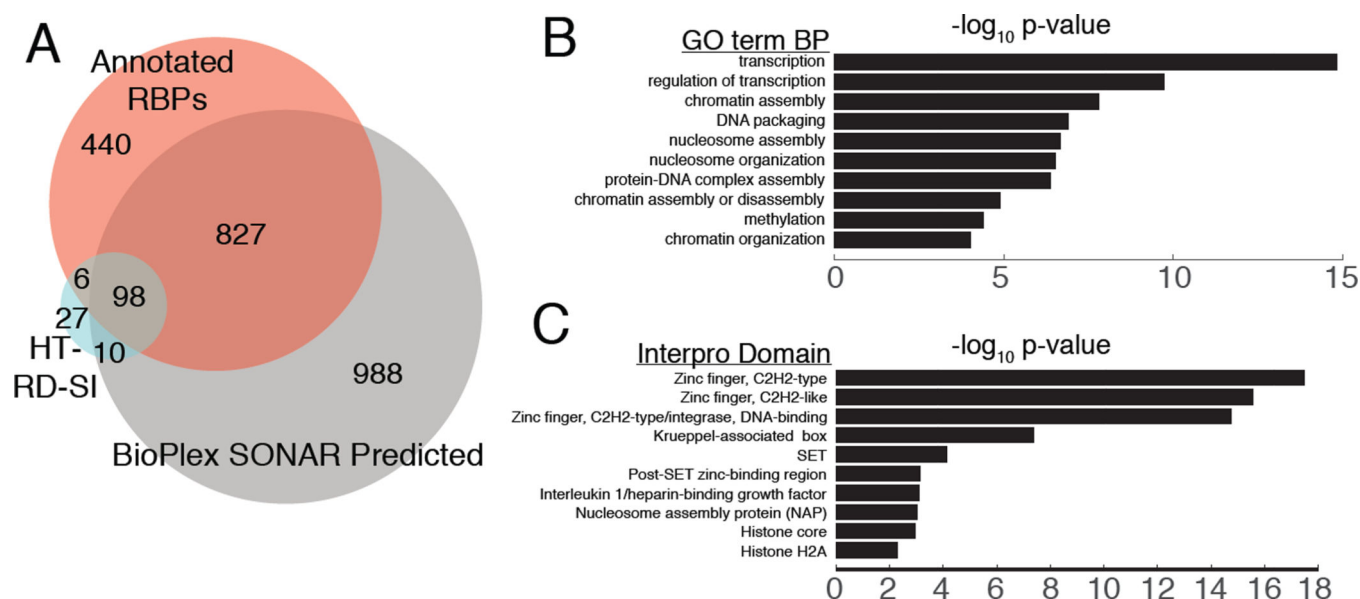


Figure 5. SONAR predicts human candidate RBPs enriched for proteins with zinc finger and DNA binding domains

A. Venn diagram showing overlap between all HT-RBP RNA-dependent SI proteins contained in the BioPlex network (light blue), all annotated RBPs contained in the BioPlex network (red), and all BioPlex SONAR predicted RBPs (at RCS>0.79; grey).

B. Bar graph displaying $-\log_{10}$ p-values for GO Biological Process (BP) terms enriched in the set of RBP candidates (RCS>0.79 and not previously annotated as RBPs) compared to all interactors within the BioPlex dataset.

C. Bar graph displaying $-\log_{10}$ p-values for INTERPRO protein domains enriched in the set of RBP candidates (RCS>0.79 and not previously annotated as RBPs) compared to all interactors within the BioPlex dataset.

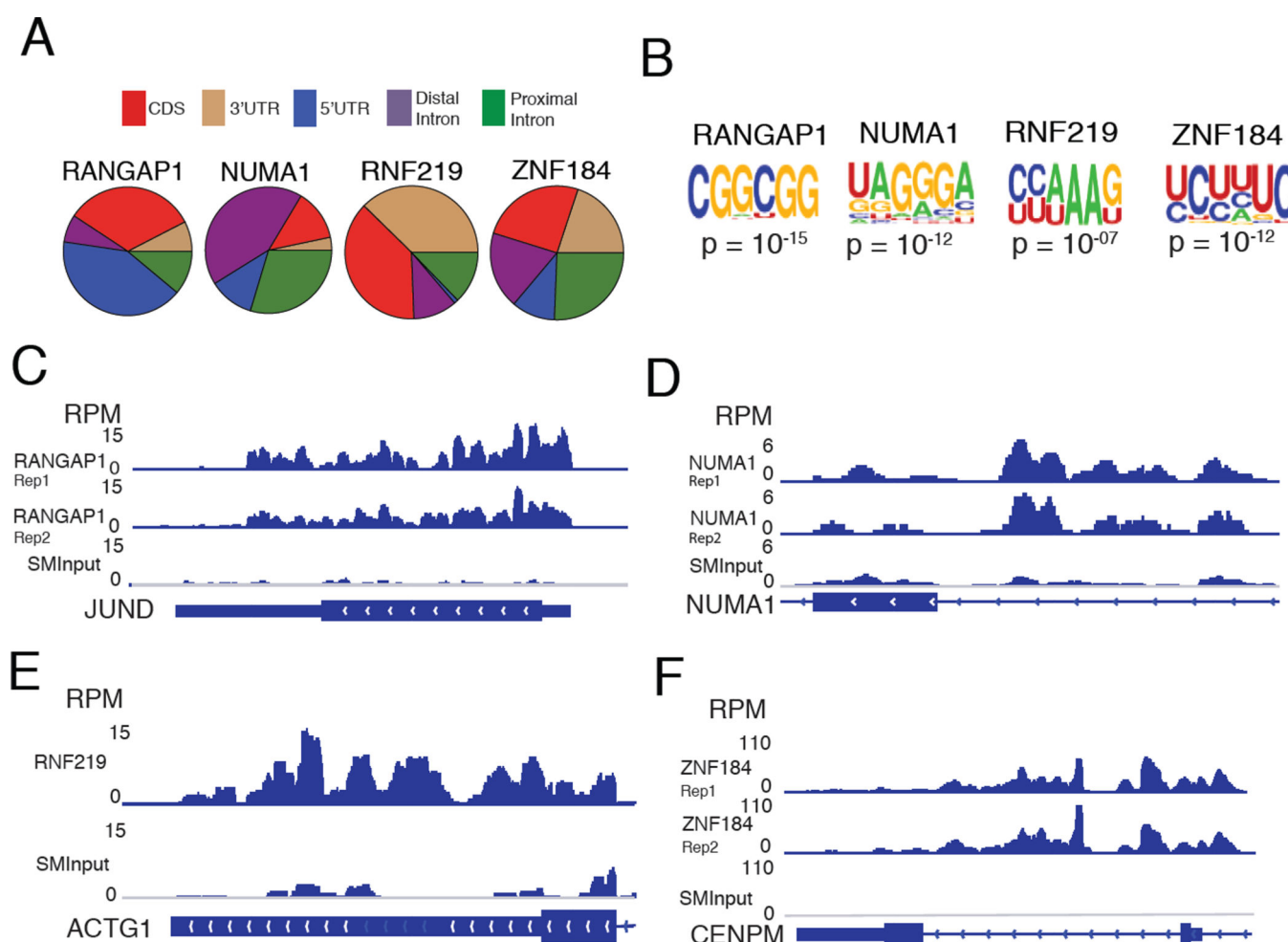


Figure 6. Enhanced CLIP validation of candidate RBPs predicted by HT-RBP interactome and SONAR classification

A. Distributions across transcript regions for peaks enriched >8 fold over size-matched input ($-\log_{10} p > 5$) from eCLIP experiments for 4 RBP candidates.

B. Motifs called and p-values for input normalized peaks described in Figure 5D.

C. Genome browser track view of RANGAP1 eCLIP data in reads per million (RPM) showing enrichment above input on the intronless JUND gene locus.

D. Genome browser track view of NUMA1 eCLIP data in reads per million (RPM) showing enrichment above input on a NUMA1 intron.

E. Genome browser track view of RNF219 eCLIP data in reads per million (RPM) showing enrichment above input on the ACTG1 3'UTR region.

F. Genome browser track view of ZNF184 eCLIP data in reads per million (RPM) showing enrichment above input on the CENPM distal intron.