

# EuRBPDB: a comprehensive resource for annotation, functional and oncological investigation of eukaryotic RNA binding proteins (RBPs)

Jian-You Liao<sup>1,2,\*</sup>, Bing Yang<sup>1,2,†</sup>, Yu-Chan Zhang<sup>3,†</sup>, Xiao-Juan Wang<sup>1,2</sup>, Yushan Ye<sup>1,4</sup>, Jing-Wen Peng<sup>1,2</sup>, Zhi-Zhi Yang<sup>1,2</sup>, Jie-Hua He<sup>1,2</sup>, Yin Zhang<sup>1,2</sup>, KaiShun Hu<sup>1,2</sup>, De-Chen Lin<sup>5,\*</sup> and Dong Yin<sup>1,2,\*</sup>

<sup>1</sup>Guangdong Provincial Key Laboratory of Malignant Tumor Epigenetics and Gene Regulation, Sun Yat-Sen Memorial Hospital, Sun Yat-Sen University, Guangzhou 510120, China, <sup>2</sup>Medical Research Center, Sun Yat-Sen Memorial Hospital, Sun Yat-Sen University, Guangzhou 510120, China, <sup>3</sup>State Key Laboratory for Biocontrol, School of Life Science, Sun Yat-Sen University, Guangzhou 510275, China, <sup>4</sup>Department of stomatology, Sun Yat-Sen Memorial Hospital, Sun Yat-sen University, Guangzhou 510120, China and <sup>5</sup>Department of Medicine, Cedars-Sinai Medical Center, Los Angeles, CA, USA

Received July 18, 2019; Revised September 05, 2019; Editorial Decision September 09, 2019; Accepted October 06, 2019

## ABSTRACT

RNA binding proteins (RBPs) are a large protein family that plays important roles at almost all levels of gene regulation through interacting with RNAs, and contributes to numerous biological processes. However, the complete list of eukaryotic RBPs including human is still unavailable. Here, we systematically identified RBPs in 162 eukaryotic species based on both computational analysis of RNA binding domains (RBDs) and large-scale RNA binding proteomic data, and established a comprehensive eukaryotic RBP database, EuRBPDB (<http://EuRBPDB.syshospital.org>). We identified a total of 311 571 RBPs with RBDs (corresponding to 6368 ortholog groups) and 3,651 non-canonical RBPs without known RBDs. EuRBPDB provides detailed annotations for each RBP, including basic information and functional annotation. Moreover, we systematically investigated RBPs in the context of cancer biology based on published literatures, PPI-network and large-scale omics data. To facilitate the exploration of the clinical relevance of RBPs, we additionally designed a cancer web interface to systematically and interactively display the biological features of RBPs in various types of cancers. EuRBPDB has a user-friendly web interface with browse and search functions, as well as data downloading function. We ex-

pect that EuRBPDB will be a widely-used resource and platform for both the communities of RNA biology and cancer biology.

## INTRODUCTION

RNA binding proteins (RBPs) are involved in the regulation of the metabolism, transportation, translation and function of both coding and non-coding RNAs through direct RNA-protein interaction (1). RBPs ensure the smooth flowing of genetic information from DNA to RNA, and ultimately to proteins, making them essential and instrumental for all physiological and pathological processes (1). Numerous diseases have been caused by the aberrant of expression or function of RBPs, including cancer, metabolic disorders and neuropathies (2–4).

Comprehensive identification and annotation of all RBPs are primary and crucial steps for characterization of their functions. To date, several RBPs databases exist for a few eukaryotes, but these databases only collected a small number of well-characterized RBPs from one or few species. For example, RBPDB is a database focusing on the collection of experimentally validated RBPs and RNA binding domains (RBDs), and it contained only 1171 RBPs from human, mouse, fly and worm (5). ATTRACT is a manually curated database that collects compiled information for only 370 well-characterized RBPs from 39 species (6). Clearly, the RBP repertoire collected by these existing databases are far from complete for any species, human included.

\*To whom correspondence should be addressed. Tel: +86 18922182515; Email: yind3@mail.sysu.edu.cn  
Correspondence may also be addressed to De-Chen Lin. Tel: +1 310 423 7740; Email: dchlin11@gmail.com  
Correspondence may also be addressed to Jian-You Liao. Tel: +86 1358054805; Email: liaojy3@mail.sysu.edu.cn  
†The authors wish it to be known that, in their opinion, the first three authors should be regarded as Joint First Authors.

RBPs bind to RNA via structurally well-defined RBDs, such as Dead box helicase domain, RNA recognition motif (RRM) (7,8). Here, we annotated proteins containing a RBD as canonical RBPs. Additionally, many studies have suggested the existence of complex protein-RNA interactions that do not require canonical RBDs (9,10), instead through other structures such as intrinsically disordered regions (IDRs) (11). It is thus challenging to identify non-canonical RBPs without known RBDs in a high-throughput and unbiased manner. Recent advances in RNA binding proteome (RBPome) technology significantly facilitate the large-scale identification of non-canonical RBPs (12–18), including the capture of polyadenylated RNA interactome (11,16–21), click chemistry-based capture of RNA interactome (13), and orthogonal organic phase separation (OOPS) of RBPs (14,15,19). These methods crosslink the RBPs with RNA using UV, then apply different strategies to extract total RBPs from cells or tissues. The purified total RBPs are used to analyze the RBPome based on mass spectrometry (MS). These RBPome technologies have been applied to many eukaryotes, including human (11,15,16,19–21), mouse (12) and fly (18), and identified a large number of novel canonical and non-canonical RBPs. It should be noted that as an experimental method, none of RBPome technologies is capable of capturing the complete category of RBPs, due to the limitation of total RBP purification strategy and MS technology (12–19). Moreover, most of the present RBPome studies applied stringent filtering process to control for the false positivity, which is associated with high false negativity and low sensitivity.

In the rapid progression of RNA biology field (1), a great need exists to build a comprehensive eukaryotic RBP database to explore the annotation, expression and function of RBPs. To address this, we collected a full list of RBDs from both Pfam (22) and published RBPome datasets from 6 eukaryotes (human, mouse, zebrafish, yeast, fly and worm) (Supplemental Table S1). In parallel, we predicted RBPs based on RBDs using HMMER (23) from the genomes of 162 eukaryotes. Upon integration, we established currently the most comprehensive database of eukaryotic RBP, EuRBPDB (Figure 1). EuRBPDB contains a total of 315 222 RBPs, with detailed annotations for each RBP. Moreover, given the crucial role of RBP in cancer biology, in order to facilitate users to explore the clinical relevance of RBPs, we separately built a Cancer web interface to display integrated cancer-associated omics datasets. The database has a user-friendly interface to interactively exhibit and search the detailed annotations. EuRBPDB will therefore greatly promote the investigation and understanding of the RNA biology.

## MATERIALS AND METHODS

### Identification and annotation of RBPs

All protein sequences of 162 eukaryotes were downloaded from Ensembl database (24) (release 96, <http://www.ensembl.org/>). Proteins were annotated as canonical RBPs if they contain one or more domains known to directly interact with RNA. The search of RBPs was based on the searching of sequence homologs of known RBDs in

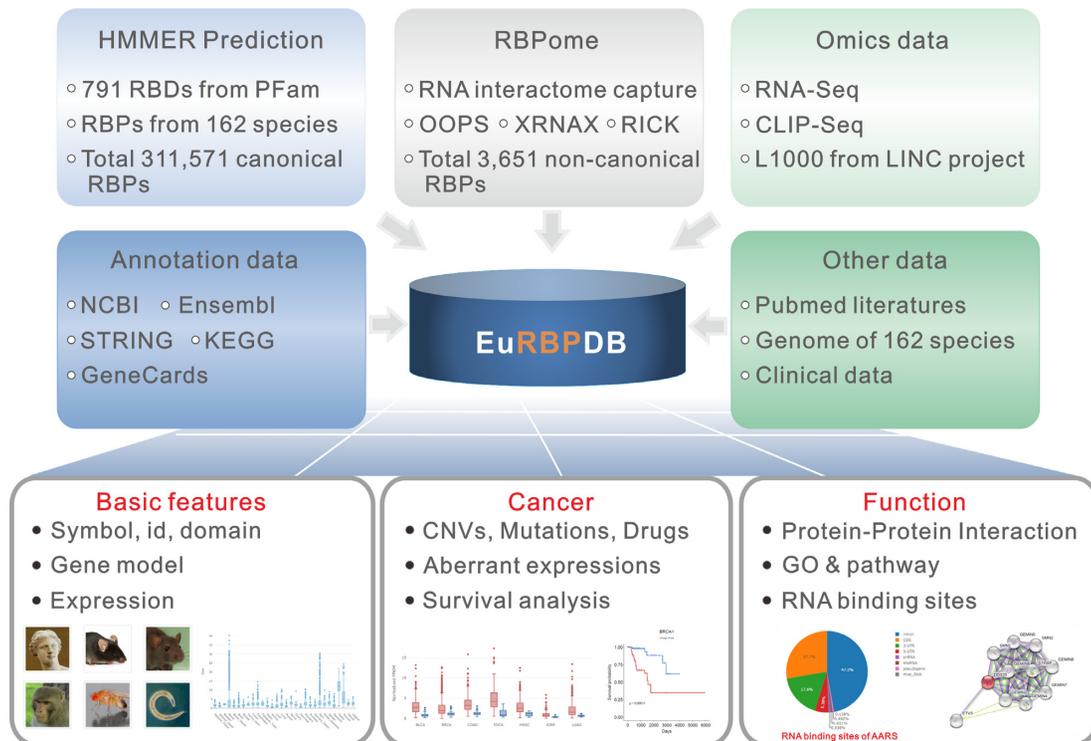
proteins using probabilistic models known as profile hidden Markov models [4]. The present RBD list was curated based on the comprehensive RBD list established by Gerstberger *et al.* (25). After careful examination, we found that eight RBDs (RRM6, KH\_3, MRL1, Ribosomal\_S3\_N, Lactamase\_B2, tRNA\_synt\_2b, RnaseH, tRNA\_anti) have been removed by Pfam, and thus they were eliminated from our list. Finally, we obtained a total of 791 RBDs (can be downloaded from [http://EuRBPDB.syshospital.org/data/download/791\\_RBDs.Pfam.gz](http://EuRBPDB.syshospital.org/data/download/791_RBDs.Pfam.gz)). We extracted RBD HMM profiles from the Protein families (Pfam) database (Pfam HMM profiles, release v32) (22), and applied the hmmsearch program in HMMER (v3.2.1) (23) package to search for all of the eukaryotic protein sequences against the RBD HMM profiles to identify RBPs. Proteins with *E*-value less than 0.0001 were considered as bona fide canonical RBPs. In total, we identified 311 571 canonical RBPs from 162 eukaryotic species. In parallel, we manually collected large-scale RBPome datasets of human, mouse, zebrafish, yeast, fly and worm from 21 published works (Supplementary Table S1). Human non-canonical RBPs are required to be detected in at least two RBPome datasets. For other species, the RBPs detected by any RBPome were included in EuRBPDB. As a result, we obtained 3651 non-canonical RBPs from six species. Finally, EuRBPDB collected a total of 315 222 RBPs, representing the largest eukaryotic RBP database currently available. EuRBPDB has four lines of evidence of RNA-binding for each RBP, namely (i) literatures supporting of RNA-binding capacity, (ii) RNA-binding domain, (iii) RBPome and (iv) RNA-binding sites detected by CLIP-Seq. We graded those RBPs with only one of four pieces of evidence as ‘putative’ in Description section on the Basic information subpage. The basic information, GO and phenotype annotation of RBPs were obtained from NCBI, Genecards and Ensembl databases. The protein–protein interaction (PPI) information was parsed from STRING database (26). The pathway annotation was obtained from KEGG database (27). Expression data were obtained from GTEx (28) and SRA.

### Classification of eukaryotic RBP family

We characterized and classified canonical RBPs by their sequence-specific RBDs. RBP family was named as the RBD domain if its RBPs only contain one type of RBD. If a RBP contains multiple types of RBDs, it was categorized into each of the family. All non-canonical RBPs were classified as non-canonical RBP family. In total, we obtained 686 RBP families.

### Orthologs and paralogs

The reciprocal best hit (RBH) method (29) was used to predict the putative orthologs of RBPs among different species. We performed the all-against-all BLASTP (v2.7.1+) search between proteins of two genomes with strict cutoffs (*E*-value  $\leq 1e-6$ , coverage  $\geq 50\%$ , identity  $\geq 30\%$ ) and annotated the reciprocal best hit pairs as orthologs. Paralogs was predicted by the BLAST score ratio (BSR) (30) approach. BLASTP search was conducted in each genome with the same parameters as in orthologs search. The BSR value cutoff was set to 0.4.



**Figure 1.** A system-level overview of the EuRBPDB core framework. A total of 315 222 RBPs, including 311 571 canonical RBPs and 3651 non-canonical RBPs, were identified by combination of computational RBP searching with RBPome profiling. All RBPs were annotated by information retrieved from public database, like NCBI, Ensembl, STRING, KEGG and GeneCards. Cancer-relevant RBPs were identified by literature mining and systematic TCGA data analysis. All the results generated by EuRBPDB were deposited in MySQL relational databases and displayed in the web pages. All species photos were downloaded from Ensembl database (24).

### Differential expression, copy number variation (CNV), mutation and survival analysis of RBPs

RNA-Seq, whole-exon sequencing and clinical data were retrieved from TCGA database using R/Bioconductor package TCGAAbiolinks (v2.8.4) (31). Differential expression analysis was performed using R package edgeR (v3.22.5) (32) [false discovery rate (FDR)  $\leq 1e-5$ ,  $\log_2$  fold change ( $\log_2$ FC)  $> = 1$ ]. Kaplan–Meier survival analysis was performed by R package survival (v2.43-3). Significant amplification and deletion genomic regions in cancer samples were downloaded from Broad GDAC Firehose website (<https://gdac.broadinstitute.org/>).

### Cellular effects of drugs to RBP expression

Two L1000 assay level-5 datasets (GSE92742 and GSE70138) (33) generated by the Library of Integrated Cellular Signatures (LINCS) project were downloaded from GEO. These datasets contain over 1 600 000 sub-datasets measuring the effects 30 744 drugs on the RNA profiles of 44 cell lines. L1000 assay datasets were parsed and displayed by campR (v1.0.1) and ggplot2 (v3.1.0) R packages as suggested by LINCS project. Expression of RBPs is displayed as  $z$ -score.

### RNA binding sites of RBPs

A total of 227 eCLIP isogenic replicated datasets generated from K562 (120 RBPs) and HepG2 (103 RBPs)

cell lines and human adrenal gland tissues (two RBPs) were retrieved from ENCODE database (<https://www.encodeproject.org/>). Peak and bam files of each datasets were downloaded. We used intersectBed of bedtools package (v2.27.1) (34) to annotate each peak, and used coverageBed of bedtools to retrieve the RPM value of each peak.

### Literature analysis of RBP

Literature mining was conducted in geneclip3 (<http://ci.smu.edu.cn/geneclip3/>). In brief, Entrez ids of all RBPs were submitted to geneclip3. Key words of function model of geneclip3 were set as ‘cancer or tumor’ to search for cancer-associated literatures, and ‘RNA binding or RNA-binding’ to search for literatures on RNA-binding. Geneclip3 was run in GeneRIF mode to search for cancer-associated literatures, and in MEDLINE mode to search for literatures on RNA-binding. The searching will return the PubMed IDs of all literatures that study the RBPs in cancers or RNA-binding capacity. The information of all literatures was retrieved from PubMed based on PubMed IDs. RBPs reported in 3 cancer-relevant studies were considered to be cancer-associated.

## DATABASE CONTENT AND WEB INTERFACE

### The web-based exploration of RBPs

EuRBPDB provides genome-wide identification of RBPs in large amount of eukaryotic species based on HMMER

searching results combined with RBPome datasets analyses. In total, 315 222 RBPs, including 311 571 canonical RBPs corresponding to 6368 ortholog groups and 3651 non-canonical RBPs, were identified in 162 eukaryotic species. With the systematic annotation of these RBPs, we designed a user-friendly web interface for users to query the database conveniently and interactively. Users can either browse the entire RBP list of any 162 eukaryotes collected in database, or search for any RBP in any eukaryotes of interest. EuRBPDB provides two different ways to browse the data, one is to browse by species, the other is to browse by family defined by RBDs. On the ‘Species’ page, 162 species were classified into 12 categories according to Ensembl taxonomy. To browse the RBP list of each species, users just need to click the species image of interest, and retrieve the detailed RBP information through the following steps: families → family gene list → single gene annotation. On the ‘Family’ page, EuRBPDB lists all 686 RBP families from 162 eukaryotes. RBP families were ordered by family size in descending order. By clicking the family name, users will get all RBPs grouped by species in this family. Users can also obtain the detailed information of RBP through the following steps: species → gene list → single gene annotation.

Users can search the specific RBP of interest using the quick search box at the top right corner of navigation bar in any page, the search will return all RBPs in any species matching the searching criteria. To browse the detailed information of any specific RBP, users can specify both the species and RBP name/ID in ‘Search’ page. Both search and browser functions direct users to the detailed information page of any specific RBP. This page comprises of two subpages, namely ‘Basic Information’ subpage and ‘Cancer Related Information’ subpage (only for human RBPs currently). All two subpages consist of a number of information sections constructed by data collected from other published databases. We can readily add any new sections to these subpages, and thus it is easy and convenient to update EuRBPDB regularly. In Basic information subpage, EuRBPDB provides basic information including gene structure (Gene Model section), evidences for RNA-binding (RBDs, RBPome, RPI and Literatures sections), expression (Expression section), and functional annotation (PPI, Pathway and Gene Ontology sections etc.). ‘Cancer Related Information’ subpage will be introduced in the following sections.

### Cancer web interface

RBPs contribute extensively and significantly to numerous processes in cancer biology. To facilitate RBP research in cancer, EuRBPDB provides cancer associated annotation of RBPs in Cancer web interface. Through systematic literature mining using geneclip3 (<http://ci.smu.edu.cn/geneclip3/>), we found that a total of 727 RBPs are reported to be associated with human cancers (reported by at least three literatures). Among them, 144 RBPs were frequently investigated (reported by >20 literatures). Moreover, we conducted differential expression, somatic mutation, CNV, as well as survival analysis based on TCGA data to reveal comprehensively the alterations of RBPs in human cancers. As a result, we identified 1361 RBPs showing aberrant expres-

sion in at least one cancer type, 2900 RBPs harboring non-sense and/or missense mutations (1761 of them mutated in RBD regions), 2851 RBPs having genomic deletions or amplifications, and 2897 RBPs exhibiting significant survival correlation.

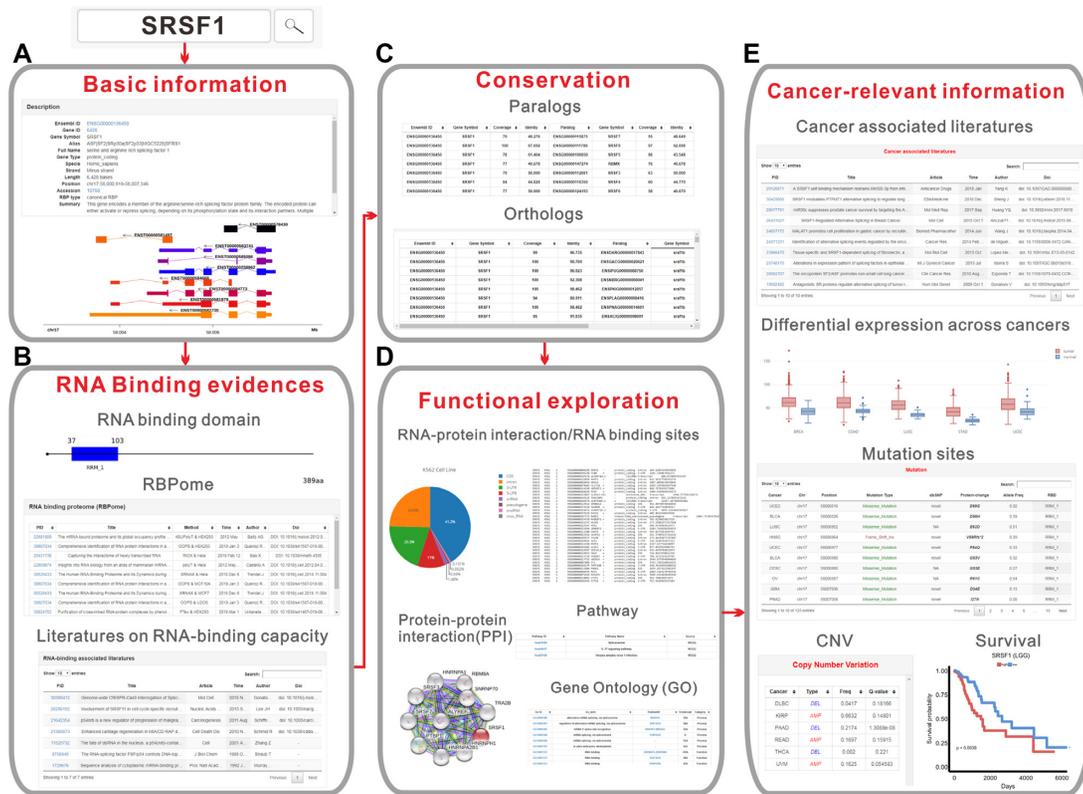
Mutational analysis of RBDs showed that certain cancer types such as Pheochromocytoma and Paraganglioma (PCGP) and PAAD, have higher mutational rate targeting RBD regions than others (Supplementary Figure S1A). This result is congruent with the findings that the expression and functions of RBPs have cell-type specificity (12,18). On the other hand, certain RBDs have higher mutation rates across human cancers (Supplementary Figure S1B), such as MMR\_HSR1 and RRM\_1 domain. Notably, mutations in RRM\_1 of RBM10 have been suggested to play important role in the development and progression of lung adenocarcinomas (35–37), highlighting that our analysis is capable of identifying functional mutations in cancer-associated RBPs. These results together suggest further investigation of the functional significance of candidate RBPs and RBD in cancer biology.

It is conceivable that larger number of genes mutated in a given RBP PPI network will result in higher degree of network dysregulation. A bar plot showing the number of aberrant RBP PPI network (defined as number of mutated gene >30% within the network) of each cancer is provided in Cancer interface. To facilitate the users to explore the number of mutated genes of each RBP PPI network in each cancer type, we added a bar plot under the PPI network figure in Basic information subpage of each RBP.

Among RBPs with cancer-associated alterations, most of them have hitherto not been reported to be associated with any cancers, providing a valuable and novel resource for cancer researchers. EuRBPDB provides the overview of the cancer-associated RBPs in ‘Cancer’ page, as well as the list of published and novel cancer-associated RBPs deposited in EuRBPDB. By clicking the ‘Details’ link of each RBP, users can be redirected to detailed information page of RBP with Cancer Related Information subpage. There are six sections in this subpage, showing the literatures investigating selected RBP (Literatures), differential expression boxplot (Differential Expression), mutations in RBP (mutation), copy number variation (CNV), survival analysis (survival), as well as the expression changes across 44 different cell lines under the treatment of ~2000 drugs (33).

### RBPredictor web-server for the annotation of eukaryotic RBPs

A web-based tool, RBPredictor, was further developed to assist users to determine whether the protein of interest (from any eukaryote) is a putative canonical RBP. Such RBP prediction is based on the RBD sets used in this study, and we performed hmm-search program in HMMER (v3.2.1) package to determine whether the protein sequence submitted is a putative RBP (25). In ‘RBPredictor’ page, users are only required to input one or multiple protein sequences in fasta format, or submit a fasta file with protein sequences. If an input protein is identified as a putative RBP, RBPredictor will also list all potential RBDs such protein harbors.



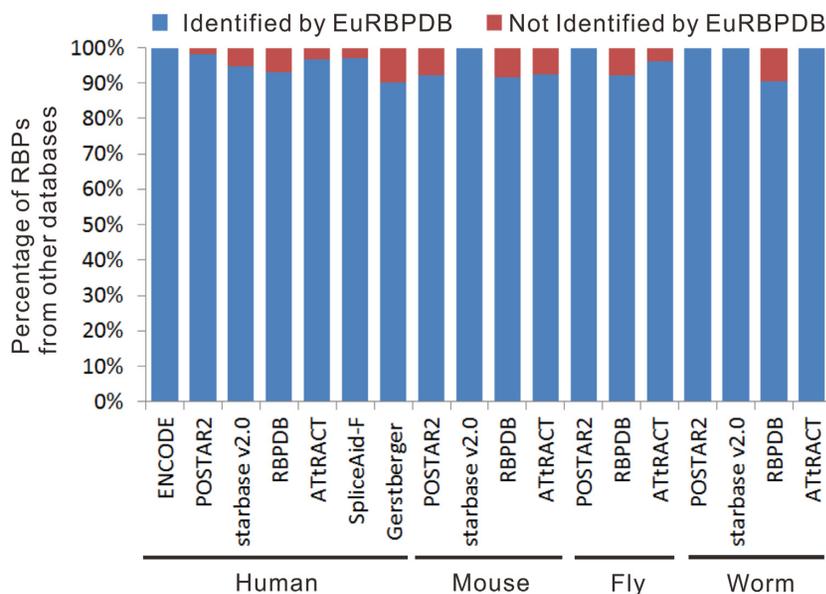
**Figure 2.** Illustration of RBP exploration in EuRBPDB: SRSF1 as an example. Through EuRBPDB searching, (A) first, users can easily obtain the detailed basic information of SRSF1, (B) then users can obtain multiple lines of RNA binding evidences of the SRSF1 including canonical RNA binding domain it contains, all RBPome datasets that detected SRSF1, and literatures that reported the binding RNA capacity of SRSF1. (C) Next, users can find the conservation status of SRSF1 through looking over the list of paralog and ortholog of SRSF1. (D) Users can further explore the function of SRSF1 through systematically investigate the protein-RNA interaction/RNA binding sites of RBP, protein-protein interaction network, pathway and gene ontology (GO) information provides by EuRBPDB. (E) Finally, users can systematically acquire the aberrant information of SRSF1 across cancers, including differential expression, mutations, copy number variations, survival correlation and literatures that reported the cancer regulatory role of SRSF1.

**DISCUSSION AND CONCLUSIONS**

In this study, we systematically identified eukaryotic RBPs by integrating both large-scale RBPome experimental data and computational RBD identification data. We identified a total of 311 571 high-confidence canonical RBPs corresponding to 6368 ortholog groups in 162 eukaryotes, and 3651 non-canonical RBPs without known RBDs in six eukaryotes (human, mouse, zebrafish, fly, worm and yeast). Currently, all non-canonical RBPs were grouped into non-canonical\_RBP protein family. 311 571 canonical RBPs formed 686 protein families. Except some large RBP families, such as RRM.1 (33 193 RBPs, 597 ortholog groups), zf-met (20 101 RBPs, 589 ortholog groups), zf-C2H2 (16 879 RBPs, 507 ortholog groups), MMR\_HSR1 (22 986 RBPs, 467 ortholog groups), most RBP families contain small amount of ortholog group (median: 4) (Supplementary Figure S2). 2961 RBPs were identified in human with high confidence, including 1836 canonical RBPs and 1135 non-canonical RBPs, significantly expanding the human RBP repertoire. Moreover, most human RBPs were found to have cancer-related alterations. We systematically annotated all eukaryotic RBPs in this study, and constructed the most comprehensive eukaryotic RBPs database, EuRBPDB. Through the integration of various large-scale

omics data (such as CLIP-Seq, RNA-Seq and L1000 assay), EuRBPDB provides a comprehensive platform to explore the function and cancer-relevance of RBPs. Users can readily obtain basic, functional and cancer-relevant information of any RBPs of interest from EuRBPDB (Figure 2). EuRBPDB also provides a RBPredictor web-server, which enables users to easily and rapidly determine whether a eukaryote protein not included in EuRBPDB is an RBP. EuRBPDB provides a framework to systematically identify eukaryotic RBPs based on RBD searching and RBPome data.

Identification of RBP through RBD matching is a highly effective and accurate approach (25). However, recent RBPome studies showed that a large number of proteins without canonical RBDs also bind RNA, and many of them bind RNA through IDRs (11). Therefore, clearly it is insufficient to identify RBPs merely based on RBD searching. On the other hand, RBPome methods are likewise incapable of detecting all RBPs because of the (i) context-dependent RNA binding capacity of many RBP approach (1); (ii) restricted expression pattern of RBPs, since the RBPome were performed in only a few cell types; (iii) technical limitation of purification strategy of total RBP (14,15,19); (iv) low sensitivity of MS technology. We also find that only about half of human canonical RBPs can be detected by different RBPome methods (Supplementary Figure S3, Supple-



**Figure 3.** EuRBPDB contained most of RBPs deposited in other databases. 100% (157/157), 98.3% (168/171), 94.7% (36/38), 93.0% (385/414), 96.9% (154/159), 97.0% (65/67), and 90.1% (1389/1542) human RBPs from ENCODE database, POSTAR2, starBase v2.0, RBPDB, ATtRACT, SpliceAid-F and Gerstberger *et al.* RBP sets (25) were contained, respectively, in EuRBPDB. 92.31% (36/39), 100% (14/14), 91.7% (373/407) and 92.6% (25/27) mouse RBPs from POSTAR2, starBase v2.0, RBPDB and ATtRACT were included, respectively, in EuRBPDB. 100% (3/3), 92.2% (226/245) and 96.2% (51/53) *Drosophila melanogaster* RBPs from POSTAR2, RBPDB and ATtRACT were included, respectively, in EuRBPDB. 100% (5/5), 100% (2/2), 90.4% (208/230) and 100% (20/20) *Caenorhabditis elegans* RBPs from POSTAR2, starBase v2.0, RBPDB and ATtRACT were included, respectively, in EuRBPDB.

mentary Table S1). Thus, presently a comprehensive way to acquire a more complete RBP repertoire is to combine the computational RBP searching with RBPome profiling.

To verify the reliability of RBP dataset we generated, we have cross-checked against all current RBP databases. The results showed that EuRBPDB identified the vast majority of the RBPs (ranging from 90.1% to 100%) across different species collected by other databases (Figure 3), validating the accuracy and consistency of our work. Furthermore, we used the GO annotation to evaluate the robustness and accuracy of our human RBP list. Indeed, we found that 95.3% of canonical RBP and 73.8% of non-canonical RBP were annotated by RNA-related GO terms, such as ‘RNA-binding’, ‘RNA modification’ and ‘endoribonuclease activity’. These results together highlight that our RBP identification approach has high accuracy and robust performance.

Many databases have been established to aid the research of RNA biology (5,6,38). However, currently no comprehensive RBP database is available for all species. All existing RBP databases focus on the collection and integration of the structure, RBD, RBP binding sites or disease correlation of small amount of well-characterized RBPs in a limited types of eukaryotes, such as RBPDB (5), ATtRACT (6), SpliceAid-F (38), POSTAR2 (39), starBase (40) etc. Compared with these RBP databases, EuRBPDB provides the largest eukaryotic RBP repertoire (315 222 RBPs, forms 6368 ortholog groups), the most comprehensive functional and cancer-associated annotation, and an intuitive and easy-to-use web interface. Therefore, EuRBPDB provides a powerful platform to decode the RBP function and regulatory mechanisms.

## FUTURE DIRECTIONS

EuRBPDB is a comprehensive eukaryotic RBP database, characterizing RBPs of 162 eukaryotic genome-wide. With the ever-increasing amount of RBPome and eukaryotic genome data, we will continue to update and maintain the RBP repertoire and annotation regularly. We will also integrate additional omics datasets (e.g. CLIP-seq, RNA-Seq) from public databases like Gene Expression Omnibus (GEO) and Sequence Read Archive (SRA) to further improve our understanding of the function and regulatory mechanism of RBPs.

## DATA AVAILABILITY

EuRBPDB database is freely available at <http://EuRBPDB.syshospital.org>.

## SUPPLEMENTARY DATA

Supplementary Data are available at NAR Online.

## FUNDING

Natural Science Foundation of China [81872140, 81420108026, 81572484, 81621004 to D.Y., 81872155, 81672621 to J.Y.L., 31770883 to Y.C.Z., 81503651 to Y.Y.]; Guangzhou Bureau of Science and Information Technology [201704030036 to D.Y., 201710010029 to Y.C.Z.]; Guangdong Science and Technology Department [2019B020226003 to D.Y., 2017B030314026]; Tip-top Scientific and Technical Innovative Youth Talents of Guangdong special support program [2016TQ03R686 to

J.Y.L., 2017TQ04N79 to Y.C.Z.]. Funding for open access charge: Natural Science Foundation of China.  
*Conflict of interest statement.* None declared.

## REFERENCES

- Hentze, M.W., Castello, A., Schwarzl, T. and Preiss, T. (2018) A brave new world of RNA-binding proteins. *Nat. Rev. Mol. Cell Biol.*, **19**, 327–341.
- Pereira, B., Billaud, M. and Almeida, R. (2017) RNA-Binding proteins in Cancer: Old players and new actors. *Trends Cancer*, **3**, 506–528.
- Cooper, T.A., Wan, L. and Dreyfuss, G. (2009) RNA and disease. *Cell*, **136**, 777–793.
- Lukong, K.E., Chang, K.W., Khandjian, E.W. and Richard, S. (2008) RNA-binding proteins in human genetic disease. *Trends Genet.*, **24**, 416–425.
- Cook, K.B., Kazan, H., Zuberi, K., Morris, Q. and Hughes, T.R. (2011) RBPDB: a database of RNA-binding specificities. *Nucleic Acids Res.*, **39**, D301–D308.
- Giudice, G., Sanchez-Cabo, F., Torroja, C. and Lara-Pezzi, E. (2016) ATtRACT-a database of RNA-binding proteins and associated motifs. *Database*, **2016**, baw035.
- Clery, A., Blatter, M. and Allain, F.H. (2008) RNA recognition motifs: boring? Not quite. *Curr. Opin. Struct. Biol.*, **18**, 290–298.
- Linder, P. and Jankowsky, E. (2011) From unwinding to clamping - the DEAD box RNA helicase family. *Nat. Rev. Mol. Cell Biol.*, **12**, 505–516.
- Ramakrishnan, V. (2014) The ribosome emerges from a black box. *Cell*, **159**, 979–984.
- Papasaikas, P. and Valcarcel, J. (2016) The Spliceosome: the ultimate RNA chaperone and sculptor. *Trends Biochem. Sci.*, **41**, 33–45.
- Castello, A., Fischer, B., Frese, C.K., Horos, R., Alleaume, A.M., Foehr, S., Curk, T., Krijgsveld, J. and Hentze, M.W. (2016) Comprehensive identification of RNA-binding domains in human cells. *Mol. Cell*, **63**, 696–710.
- Kwon, S.C., Yi, H., Eichelbaum, K., Fohr, S., Fischer, B., You, K.T., Castello, A., Krijgsveld, J., Hentze, M.W. and Kim, V.N. (2013) The RNA-binding protein repertoire of embryonic stem cells. *Nat. Struct. Mol. Biol.*, **20**, 1122–1130.
- Bao, X., Guo, X., Yin, M., Tariq, M., Lai, Y., Kanwal, S., Zhou, J., Li, N., Lv, Y., Pulido-Quetglas, C. et al. (2018) Capturing the interactome of newly transcribed RNA. *Nat. Methods*, **15**, 213–220.
- Urdaneta, E.C., Vieira-Vieira, C.H., Hick, T., Wessels, H.H., Figini, D., Moschall, R., Medenbach, J., Ohler, U., Granneman, S., Selbach, M. et al. (2019) Purification of cross-linked RNA-protein complexes by phenol-toluol extraction. *Nat. Commun.*, **10**, 990.
- Queiroz, R.M.L., Smith, T., Villanueva, E., Marti-Solano, M., Monti, M., Pizzinga, M., Mirea, D.M., Ramakrishna, M., Harvey, R.F., Dezi, V. et al. (2019) Comprehensive identification of RNA-protein interactions in any organism using orthogonal organic phase separation (OOPS). *Nat. Biotechnol.*, **37**, 169–178.
- Castello, A., Fischer, B., Eichelbaum, K., Horos, R., Beckmann, B.M., Strein, C., Davey, N.E., Humphreys, D.T., Preiss, T., Steinmetz, L.M. et al. (2012) Insights into RNA biology from an atlas of mammalian mRNA-binding proteins. *Cell*, **149**, 1393–1406.
- Baltz, A.G., Munschauer, M., Schwanhauser, B., Vasile, A., Murakawa, Y., Schueler, M., Youngs, N., Penfold-Brown, D., Drew, K., Milek, M. et al. (2012) The mRNA-bound proteome and its global occupancy profile on protein-coding transcripts. *Mol. Cell*, **46**, 674–690.
- Sysoev, V.O., Fischer, B., Frese, C.K., Gupta, I., Krijgsveld, J., Hentze, M.W., Castello, A. and Ephrussi, A. (2016) Global changes of the RNA-bound proteome during the maternal-to-zygotic transition in *Drosophila*. *Nat. Commun.*, **7**, 12128.
- Trendel, J., Schwarzl, T., Horos, R., Prakash, A., Bateman, A., Hentze, M.W. and Krijgsveld, J. (2019) The human RNA-binding proteome and its dynamics during translational arrest. *Cell*, **176**, 391–403.
- Castello, A., Frese, C.K., Fischer, B., Jarvelin, A.I., Horos, R., Alleaume, A.M., Foehr, S., Curk, T., Krijgsveld, J. and Hentze, M.W. (2017) Identification of RNA-binding domains of RNA-binding proteins in cultured cells on a system-wide scale with RBDmap. *Nat. Protoc.*, **12**, 2447–2464.
- Perez-Perri, J.I., Rogell, B., Schwarzl, T., Stein, F., Zhou, Y., Rettel, M., Brosig, A. and Hentze, M.W. (2018) Discovery of RNA-binding proteins and characterization of their dynamic responses by enhanced RNA interactome capture. *Nat. Commun.*, **9**, 4408.
- El-Gebali, S., Mistry, J., Bateman, A., Eddy, S.R., Luciani, A., Potter, S.C., Qureshi, M., Richardson, L.J., Salazar, G.A., Smart, A. et al. (2019) The Pfam protein families database in 2019. *Nucleic Acids Res.*, **47**, D427–D432.
- Eddy, S.R. (1998) Profile hidden Markov models. *Bioinformatics*, **14**, 755–763.
- Zerbino, D.R., Achuthan, P., Akanni, W., Amode, M.R., Barrell, D., Bhai, J., Billis, K., Cummins, C., Gall, A., Giron, C.G. et al. (2018) Ensembl 2018. *Nucleic Acids Res.*, **46**, D754–D761.
- Gerstberger, S., Hafner, M. and Tuschl, T. (2014) A census of human RNA-binding proteins. *Nat. Rev. Genet.*, **15**, 829–845.
- Szklarczyk, D., Franceschini, A., Wyder, S., Forslund, K., Heller, D., Huerta-Cepas, J., Simonovic, M., Roth, A., Santos, A., Tsafou, K.P. et al. (2015) STRING v10: protein-protein interaction networks, integrated over the tree of life. *Nucleic Acids Res.*, **43**, D447–D452.
- Kanehisa, M., Goto, S., Kawashima, S., Okuno, Y. and Hattori, M. (2004) The KEGG resource for deciphering the genome. *Nucleic Acids Res.*, **32**, D277–D280.
- Consortium, G.T. (2013) The Genotype-Tissue Expression (GTEx) project. *Nat. Genet.*, **45**, 580–585.
- Moreno-Hagelsieb, G. and Latimer, K. (2008) Choosing BLAST options for better detection of orthologs as reciprocal best hits. *Bioinformatics*, **24**, 319–324.
- Rasko, D.A., Myers, G.S. and Ravel, J. (2005) Visualization of comparative genomic analyses by BLAST score ratio. *BMC Bioinformatics*, **6**, 2.
- Colaprico, A., Silva, T.C., Olsen, C., Garofano, L., Cava, C., Garolini, D., Sabedot, T.S., Malta, T.M., Pagnotta, S.M., Castiglioni, I. et al. (2016) TCGAAbiolinks: an R/Bioconductor package for integrative analysis of TCGA data. *Nucleic Acids Res.*, **44**, e71.
- Robinson, M.D., McCarthy, D.J. and Smyth, G.K. (2010) edgeR: a Bioconductor package for differential expression analysis of digital gene expression data. *Bioinformatics*, **26**, 139–140.
- Subramanian, A., Narayan, R., Corsello, S.M., Peck, D.D., Natoli, T.E., Lu, X., Gould, J., Davis, J.F., Tubelli, A.A., Asiedu, J.K. et al. (2017) A next generation connectivity Map: L1000 platform and the first 1,000,000 profiles. *Cell*, **171**, 1437–1452.
- Quinlan, A.R. and Hall, I.M. (2010) BEDTools: a flexible suite of utilities for comparing genomic features. *Bioinformatics*, **26**, 841–842.
- Imielinski, M., Berger, A.H., Hammerman, P.S., Hernandez, B., Pugh, T.J., Hodis, E., Cho, J., Suh, J., Capelletti, M., Sivachenko, A. et al. (2012) Mapping the hallmarks of lung adenocarcinoma with massively parallel sequencing. *Cell*, **150**, 1107–1120.
- Sebestyen, E., Singh, B., Minana, B., Pages, A., Mateo, F., Pujana, M.A., Valcarcel, J. and Eyra, E. (2016) Large-scale analysis of genome and transcriptome alterations in multiple tumors unveils novel cancer-relevant splicing networks. *Genome Res.*, **26**, 732–744.
- Hernandez, J., Bechara, E., Schlesinger, D., Delgado, J., Serrano, L. and Valcarcel, J. (2016) Tumor suppressor properties of the splicing regulatory factor RBM10. *RNA Biol.*, **13**, 466–472.
- Giulietti, M., Piva, F., D'Antonio, M., D'Onorio De Meo, P., Paoletti, D., Castrignano, T., D'Erchia, A.M., Picardi, E., Zambelli, F. and Principato, G. (2013) SpliceAid-F: a database of human splicing factors and their RNA-binding sites. *Nucleic Acids Res.*, **41**, D125–D131.
- Zhu, Y., Xu, G., Yang, Y.T., Xu, Z., Chen, X., Shi, B., Xie, D., Lu, Z.J. and Wang, P. (2019) POSTAR2: deciphering the post-transcriptional regulatory logics. *Nucleic Acids Res.*, **47**, D203–D211.
- Li, J.H., Liu, S., Zhou, H., Qu, L.H. and Yang, J.H. (2014) starBase v2.0: decoding miRNA-ceRNA, miRNA-ncRNA and protein-RNA interaction networks from large-scale CLIP-Seq data. *Nucleic Acids Res.*, **42**, D92–D97.