1 **Binding specificities of human RNA binding proteins towards structured and linear**
2 **RNA sequences**
3

4 Arttu Jolma[1],*, Jilin Zhang[1],*, Estefania Mondragón[4],*, Teemu Kivioja[2], Yimeng Yin[1], Fangjie
5 Zhu[1], Quaid Morris[5,6,7,8], Timothy R. Hughes[5,6], L. James Maher III[4] and Jussi Taipale[1,2,3,#]
6
7 *[1]Department of Medical Biochemistry and Biophysics, Karolinska Institutet, Solna, Sweden*
8 *[2]Genome-Scale Biology Program, University of Helsinki, Helsinki, Finland*
9 *[3]Department of Biochemistry, University of Cambridge, Cambridge CB2 1GA, United Kingdom*
10 *[4]Department of Biochemistry and Molecular Biology, Mayo Clinic College of Medicine and*
11 *Science, Rochester, USA*
12 *[5]Department of Molecular Genetics, University of Toronto, Toronto, Canada*
13 *[6]Donnelly Centre, University of Toronto, Toronto, Canada*
14 *[7]Edward S Rogers Sr Department of Electrical and Computer Engineering, University of*
15 *Toronto, Toronto, Canada*
16 *[8]Department of Computer Science, University of Toronto, Toronto, Canada*
17
18 *\*Authors contributed equally*
19 *#Corresponding author*

20
21
22 **ABSTRACT**
23

24 **Sequence specific RNA-binding proteins (RBPs) control many important**
25 **processes affecting gene expression. They regulate RNA metabolism at multiple levels,**
26 **by affecting splicing of nascent transcripts, RNA folding, base modification, transport,**
27 **localization, translation and stability. Despite their central role in most aspects of RNA**
28 **metabolism and function, most RBP binding specificities remain unknown or**
29 **incompletely defined. To address this, we have assembled a genome-scale collection**
30 **of RBPs and their RNA binding domains (RBDs), and assessed their specificities using**
31 **high throughput RNA-SELEX (HTR-SELEX). Approximately 70% of RBPs for which we**
32 **obtained a motif bound to short linear sequences, whereas ~30% preferred**
33 **structured motifs folding into stem-loops. We also found that many RBPs can bind to**
34 **multiple distinctly different motifs. Analysis of the matches of the motifs on human**
35 **genomic sequences suggested novel roles for many RBPs in regulation of splicing, and**
36 **also revealed RBPs that are likely to control specific classes of transcripts. Global**
37 **analysis of the motifs also revealed an enrichment of G and U nucleotides. Masking of**
38 **G and U by proteins increases the specificity of RNA folding, as both G and U can pair**
39 **to two other RNA bases via canonical Watson-Crick or G-U base pairs. The collection**
40 **containing 145 high resolution binding specificity models for 86 RBPs is the largest**
41 **systematic resource for the analysis of human RBPs, and will greatly facilitate future**
42 **analysis of the various biological roles of this important class of proteins.**

43
44
45

**INTRODUCTION**

The abundance of protein and RNA molecules in a cell depends both on their rates of production and degradation. These rates are determined directly or indirectly by the sequence of DNA. The transcription rate of RNA and the rate of degradation of proteins is determined by DNA and protein sequences, respectively. However, most regulatory steps that control gene expression are influenced by the sequence of the RNA itself. These processes include RNA splicing, localization, stability, and translation. These processes can be affected by RNA-binding proteins (RBPs) that specifically recognize short RNA sequence elements (Glisovic et al., 2008).

RBPs can recognize their target sites using two mechanisms: they can form direct contacts to the RNA bases of an unfolded RNA chain, and/or recognise folded RNA-structures (Loughlin et al., 2009). These two recognition modes are not mutually exclusive, and the same RBP can combine both mechanisms in recognition of its target sequence. The RBPs that bind to unfolded target sequences generally bind to each base independently of each other, and their specificity can thus be well explained by a simple position weight matrix (PWM) model. However, recognition of a folded RNA-sequence leads to strong positional interdependencies between different bases due to base pairing. In addition to the canonical Watson-Crick base pairs G:C and A:U, double-stranded RNA commonly contains also G:U base pairs, and can also accommodate other non-canonical base pairing configurations in specific structural contexts (Varani and McClain, 2000).

It has been estimated that the human genome encodes approximately 1500 proteins that can associate with RNA (Gerstberger et al., 2014). Only some of the RBPs are thought to be sequence specific. Many RNA-binding proteins bind only a single RNA species (e.g. ribosomal proteins), or serve a structural role in ribonucleoprotein complexes or the spliceosome. As RNA can fold to complex three-dimensional structures, defining what constitutes an RBP is not simple. In this work, we have focused on RBPs that are likely to bind to short sequence elements analogously to sequence-specific DNA binding transcription factors. The number of such RBPs can be estimated based on the number of proteins containing one or more canonical RNA-binding protein domains. The total number is likely to be ~400 RBPs (Cook et al., 2011; Ray et al., 2013). The major families of RBPs contain canonical RNA-binding protein domains (RBDs) such as the RNA recognition motif (RRM), CCCH zinc finger, K homology (KH) and cold shock domain (CSD). In addition, smaller number of proteins bind RNA using La, HEXIM, PUF, THUMP, YTH, SAM and TRIM-NHL domains. In addition, many "non-canonical" RBPs that do not contain any of the currently known RBDs have been reported to specifically to RNA (see, for example (Gerstberger et al., 2014)).

Various methods have been developed to determine the binding positions and specificities of RNA binding proteins. Methods that use crosslinking of RNA to proteins followed by immunoprecipitation and then massively parallel sequencing (CLIP-seq or HITS-CLIP, reviewed in (Darnell, 2010) and PAR-CLIP (Hafner et al., 2010) can determine RNA positions bound by RBPs *in vivo*, whereas other methods such as SELEX (Tuerk and Gold, 1990), RNA bind-N-seq (Lambert et al., 2015) and RNAcompete (Ray et al., 2009) can determine motifs bound by RBPs *in vitro*. Most high-resolution models derived to date have been determined using RNAcompete, where microarrays are used to generate a library of RNA-molecules containing all possible 7-base long subsequences in at least 256

92    oligonucleotides, and the desired RBP is then used to select its target sites followed by
93    detection of the bound sites using a second microarray. RNAcompete has been used to
94    analyze large numbers of RBPs from multiple species including generation of PWMs for 75
95    human RBPs (Ray et al., 2013).
96        The CISBP-RNA database (Ray et al., 2013) (Database Build 0.6) currently lists total
97    of 392 high-confidence RBPs in human, but contains high-resolution specificity models for
98    only 100 of them (Ray et al., 2013). In addition, a literature curation based database RBPDB
99    (Cook et al., 2011) contains experimental data for 133 human RBPs, but mostly contains
100   individual target- or consensus sites, and only has high resolution models for 39 RBPs. Thus,
101   despite the central importance of RBPs in fundamental cellular processes, the precise
102   sequence elements bound by most RBPs remain to be determined. To address this problem,
103   we have in this work developed high-throughput RNA SELEX (HTR-SELEX) and used it to
104   determine binding specificities of human RNA binding proteins. Our analysis suggests that
105   many RBPs prefer to bind structured RNA motifs, and can associate with several distinct
106   sequences. The distribution of motif matches in the genome indicates that many RBPs have
107   central roles in regulation of RNA metabolism and activity in cells.
108
109

3

110    **RESULTS**
111
112
113    **Identification of RNA-binding motifs using HTR-SELEX**
114
115    To identify binding specificities of human RBPs, we established a collection of
116    canonical and non-canonical full-length RBPs and RNA binding domains. The full-length
117    constructs representing 819 putative RBPs were picked from the Orfeome 3.1 and 8.1
118    collections (Lamesch et al., 2007) based on annotations of the CisBP database for
119    conventional RBPs (Ray et al., 2013) and Gerstberger et al. (Gerstberger et al., 2014) to
120    include additional unconventional RBPs. The 293 constructs designed to cover all canonical
121    RBDs within 156 human RBPs were synthesized based on Interpro defined protein domain
122    calls from ENSEMBL v76. Most RBD constructs contained all RBDs of a given protein with 15
123    amino-acids of flanking sequence (see **Table S1** for details). Constructs containing subsets
124    of RBDs were also analyzed for some very large RBPs. Taken together our clone collection
125    covered 942 distinct proteins. The RBPs were expressed in *E.coli* as fusion proteins with
126    thioredoxin, incorporating an N-terminal hexahistidine and a C-terminal SBP tag (Jolma et
127    al., 2015).
128    To identify RNA sequences that bind to the proteins, we subjected the proteins to
129    HTR-SELEX (**Figure 1A**). In HTR-SELEX, a 40 bp random DNA sequence containing a sample
130    index and 5' and 3' primer binding sequences is transcribed into RNA using T7 RNA
131    polymerase, and incubated with the individual proteins in the presence of RNase inhibitors,
132    followed by capture of the proteins using metal-affinity resin. After washing and RNA
133    recovery, a DNA primer is annealed to the RNA, followed by amplification of the bound
134    sequences using a reverse-transcription polymerase chain reaction (RT-PCR) using primers
135    that regenerate the T7 RNA polymerase promoter. The entire process is repeated up to a
136    total of four selection cycles. The amplified DNA is then sequenced, followed by identification
137    of motifs using the Autoseed pipeline (Nitta et al., 2015) modified to analyze only the
138    transcribed strand. Compared to previous methods such as RNAcompete, HTR-SELEX uses a
139    selection library with very high sequence complexity, allowing identification of long RNA
140    binding preferences.
141    The analysis resulted in generation of 145 binding specificity models for 86 RBPs.
142    Most of the results (66 RBPs) were replicated in a second HTR-SELEX experiment. The
143    success rate of our experiments was ∼ 22% for the canonical RBPs, whereas the fraction of
144    the successful non-canonical RBPs was much lower (∼ 1.3%; **Table S1**). Comparison of our
145    data with a previous dataset generated using RNAcompete (Ray et al., 2013) and to older
146    data that has been compiled in the RBPDB-database (Cook et al., 2011) revealed that the
147    specificities were generally consistent with the previous findings (**Figure S1**). HTR-SELEX
148    resulted in generation of a larger number of motifs than the previous systematic studies, and
149    revealed the specificities of 49 RBPs whose high-resolution specificity was not previously
150    known (**Figure 1B**). Median coverage per RBD family was 24 % (**Figure 1C**). Compared to
151    the motifs from previous studies, the motifs generated with HTR-SELEX were also wider, and
152    had a higher information content (**Figure S2**), most likely due to the fact that the sequences
153    are selected from a more complex library in HTR-SELEX (see also (Yin et al., 2017)). The
154    median width and information contents of the models were 10 bases and 10 bits,
155    respectively.

4

156
157 **Some RBPs bind to RNA as dimers or multimers**
158
159        Analysis of enriched sequences revealed that 31% of RBPs could bind to a site
160 containing a direct repeat of the same sequence (**Figure S3**), suggesting that some RBPs
161 were homodimers, or interacted to form a homodimer when bound to the RNA. In these
162 cases, the gap between the repeats was generally short, with a median gap of 5 nucleotides
163 (**Figure S3**). To determine whether the HTR-SELEX identified gap length preferences were
164 also observed in sites bound *in vivo*, we compared our data against existing *in vivo* data for
165 five RBPs for which high quality PAR-CLIP and HITS-CLIP derived data was available from
166 previous studies (Farazi et al., 2014; Hafner et al., 2010; Weyn-Vanhentenryck et al., 2014),
167 and found that preferred spacing identified in HTR-SELEX was in most cases also observed
168 in the *in vivo* data (**Figure S4**).
169
170
171 **Recognition of RNA structures by RBPs**
172
173        Unlike double-stranded DNA, whose structure is relatively independent of sequence,
174 RNA folds into complex, highly sequence-dependent three dimensional structures. To
175 analyze whether RBP binding depends on RNA secondary structure, we identified
176 characteristic patterns of dsRNA formation by identifying correlations between all two base
177 positions either within the motif or in its flanking regions, using a measure described in Nitta
178 et al., (Nitta et al., 2015) that determines how much the observed count of combinations of a
179 given set of two bases deviate from expected count based on independence of the positions
180 (**Figure 2A**). The vast majority of the observed deviations from the independence
181 assumption were consistent with the formation of an RNA stem-loop structure (example in
182 **Figure 2B**). In addition, we identified one RBP, LARP6, that bound to a predicted internal
183 loop embedded in a double-stranded RNA stem (**Figure 2C, Figure S5**). This binding
184 specificity is consistent with the earlier observation that LARP6 binds to stem-loops with
185 internal loops found in mRNAs encoding the collagen proteins COL1A1, COL1A2 and COL3A1
186 (Cai et al., 2010) (**Figure S5**).
187        In total, 69% (59 of 86) of RBPs recognized linear sequence motifs that did not appear
188 to have a preference for a specific RNA secondary structure. The remaining 31% (27 of 86)
189 of RBPs could bind at least one structured motif (**Figure 2D**); this group included several
190 known structure-specific RBPs, such as RC3H1, RC3H2 (Leppek et al., 2013), RBMY1E,
191 RBMY1F, RBMY1J (Skrisovska et al., 2007) and HNRNPA1 (Chen et al., 2016; Orenstein et al.,
192 2018). A total of 15 RBPs bound exclusively to structured motifs, whereas 12 RBPs could
193 bind to both structured and unstructured motifs. The median length of the stem region
194 observed in all motifs was 5 bp, and the loops were between 3 and 15 bases long, with a
195 median length of 11 (**Figure 2E**). Of the RBP families, KH and HEXIM motifs we found were
196 linear, whereas some proteins from RRM, CSD, Zinc finger and LA-domain families could bind
197 to both structured and unstructured sites (**Figure S6**).
198        To model RBP binding to stem-loop structures, we developed a simple stem-loop
199 model (SLM; **Figure 2B**). This model describes the loop as a position independent model
200 (PWM), and the stem by a nucleotide pair model where the frequency of each combination
201 of two bases at the paired positions is recorded. In addition, we developed two different

202    visualizations of the model, a T-shaped motif that describes the mononucleotide distribution
203    for the whole model, and the frequency of each set of bases at the paired positions by
204    thickness of edges between the bases (**Figure 3**), and a simple shaded PWM where the stem
205    part is indicated by a gray background where the darkness of the background indicates the
206    fraction of bases that pair with each other using Watson-Crick or G:U base pairs (**Figure 3**).
207    On average, the SLM increased the information content of the motifs by 4.2 bits (**Figure S7**).
208    As expected from the correlation structure, a more detailed analysis of the number of paired
209    bases within 10 bp from the seed sequence of MKRN1 revealed that >80% of individual
210    sequence reads had more than four paired bases, compared to ~15% for the control RBP
211    (ZRANB2) for which a structured motif was not identified.
212
213
214    **Classification of RBP motifs**

215    To analyze the motif collection globally, we developed PWM and SLM models for all
216    RBPs. To compare the motifs, we determined their similarity using SSTAT. To simplify the
217    analysis, the PWM models were used for this comparison even for factors that bound to the
218    structured motifs. We then used the dominating set method (Jolma et al., 2013) to identify
219    distinctly different motifs (**Figure S8**). Comparison of the motifs revealed that in general, the
220    specificities of evolutionarily related RBPs were similar (**Figure 5** and **Figure S8**). For the
221    largest RRM family, the 96 members were represented by 47 specificity classes, whereas the
222    smaller classes such as CCCH, KH, CSD, and HEXIM were represented by 9, 10, 6 and 1 motifs,
223    representing 17, 11, 7 and 2 different specificities, respectively (**Figure S8**).
224    Analysis of the dinucleotide content of all motifs revealed unexpected differences in
225    occurrence of distinct dinucleotides within the PWMs. The dinucleotides GG, GU, UG and UU
226    were far more common than other dinucleotides (**Figure 4G**; fold change 2.75; $p < 0.00225$;
227    t-test). This suggests that G and U bases are most commonly bound by RBPs. This effect was
228    not due to the presence of stem structures in the motifs, as the unstructured motifs were also
229    enriched in G and U. The masking of G and U bases by protein binding may assist in folding
230    of RNA to defined structures, as G and U bases have lower specificity in pairing compared to
231    C and A, due to the presence of the G:U base pair in RNA.
232    Most RBPs bound to only one motif. However, 41 RBPs could bind to multiple
233    different sites (**Figure 5**). In five cases, the differences between the primary and secondary
234    motif could be explained by a difference in spacing between the two half-sites. In 12 cases,
235    one of the motifs was structured, and the other linear. In addition, in eight RBPs the primary
236    and secondary motifs represented two different structured motifs, where the loop length or
237    the loop sequence varied (**Figure 5**). In addition, for four RBPs, we recovered more than two
238    different motifs. The most complex binding specificity we identified belonged to LARP6
239    (**Figure 5 and S9**), which could bind to multiple simple linear motifs, multiple dimeric motifs,
240    and the internal loop-structure described above.
241
242
243    **Conservation and occurrence of motif matches**
244
245    We next analyzed the enrichment of the motif occurrences in different classes of
246    human transcripts. The normalized density of motifs for each factor at both strands of DNA

6

247 was evaluated for transcription start sites (TSSs), splice donor and acceptor sites, and
248 translational start and stop positions (see **Supplementary Data S1** for full data). This
249 analysis revealed that many RBP recognition motifs were enriched at splice junctions. The
250 most enriched linear motif in splice donor sites was ZRANB2, a known regulator of
251 alternative splicing (**Figure 6A**) (Loughlin et al., 2009). Analysis of matches to structured
252 motifs revealed even stronger enrichment of motifs for ZC3H12A, B and C to splice donor
253 sites (**Figure 6A**). These results suggest a novel role for ZC3H12 proteins in regulation of
254 splicing. The motifs for both ZRANB2 and ZC3H12 protein factors were similar but not
255 identical to the canonical splice donor consensus sequence (ag | GURagu) that is recognized
256 by the spliceosome, suggesting that these proteins may act by masking a subset of splice
257 donor sites.
258      Analysis of splice acceptor sites also revealed that motifs for known components of
259 the spliceosome, such as RBM28 (Damianov et al., 2006), were enriched in introns and
260 depleted in exons. Several motifs were also enriched at the splice junction, including the
261 known regulators of splicing IGF2BP1 and ZFR (**Supplementary Data S1)** (Haque et al.,
262 2018; Huang et al., 2018). In addition, we found several motifs that mapped to the 5' of the
263 splice junction, including some known splicing factors such as QKI (Hayakawa-Yano et al.,
264 2017) and ELAVL1 (Bakheet et al., 2018), and some factors such as DAZL, CELF1 and BOLL
265 for which a role in splicing has to our knowledge not been reported (**Figure 6A** and
266 **Supplementary Data S1)** (Rosario et al., 2017; Xia et al., 2017).
267      To determine whether the identified binding motifs for RBPs are biologically
268 important, we analyzed the conservation of the motif matches in mammalian genomic
269 sequences close to splice junctions. This analysis revealed strong conservation of several
270 classes of motifs in the transcripts (**Figure 6B**), indicating that many of the genomic
271 sequences matching the motifs are under purifying selection.
272      Both matches to ZRANB2 and ZC3H12 motifs were also enriched in 5' regions of
273 transcripts, but not in anti-sense transcripts originating from promoters (**Figure 6C**),
274 suggesting that these motifs also have a role in differentiating between sense and anti-sense
275 transcripts of mRNAs.
276      To identify biological roles of the motifs, we also used Gene Ontology Enrichment
277 analysis to identify motifs that were enriched in specific types of mRNAs. This analysis
278 revealed that many RBP motifs are specifically enriched in particular classes of transcripts.
279 For example, we found that MEX3B motifs were enriched in genes involved in type I
280 interferon-mediated signaling pathway (**Figure 6D**). Taken together, our analysis indicates
281 that RBP motifs are biologically relevant, as matches to the motifs are conserved, and occur
282 specifically in genomic features and in transcripts having specific biological roles.
283
284
285
286

## DISCUSSION

In this work, we have determined the RNA-binding specificities of a large collection of human RNA-binding proteins. The tested proteins included both proteins with canonical RNA binding domains and putative RBPs identified experimentally (Gerstberger et al., 2014; Ray et al., 2013). The method used for analysis involved selection of RNA ligands from a collection of random 40 nucleotide sequences. Compared to previous analyses of RNA-binding proteins, the HTR-SELEX method allows identification of structured motifs, and motifs that are relatively high in information content. The method can identify simple sequence motifs or structured RNAs, provided that their information content is less than ~40 bits. However, due to the limit on information content, and requirement of relatively high-affinity binding, the method does not generally identify highly structured RNAs that in principle could bind to almost any protein. Consistent with this, most binding models that we could identify were for proteins containing canonical RBPs.

Motifs were identified for a total of 86 RBPs. Interestingly, a large fraction of all RBPs (47%) could bind to multiple distinctly different motifs. The fraction is much higher than that observed for double-stranded DNA binding transcription factors, suggesting that sequence recognition and/or individual binding domain arrangement on single-stranded RNA can be more flexible than on dsDNA. Analysis of the mononucleotide content of all the models also revealed a striking bias towards recognition of G and U over C and A. This may reflect the fact that formation of RNA structures is largely based on base pairing, and that G and U are less specific in their base pairings that C and A. Thus, RBPs that mask G and U bases increase the overall specificity of RNA folding in cells.

Similar to proteins, depending on sequence, single-stranded nucleic acids may fold into complex and stable structures, or remain largely disordered. Most RBPs preferred short linear RNA motifs, suggesting that they recognize RNA motifs found in unstructured or single-stranded regions. However, approximately 31% of all RBPs preferred at least one structured motif. The vast majority of the structures that they recognized were simple stem-loops, with relatively short stems, and loops of 3-15 bases. Most of the base specificity of the motifs was found in the loop region, with only one or few positions in the stem displaying specificity beyond that caused by the paired bases. This is consistent with the structure of fully-paired double-stranded RNA where base pair edge hydrogen-bonding information is largely inaccessible in the deep and narrow major groove. In addition, we identified one RBP that bound to a more complex structure. LARP6 recognized an internal loop structure where two base-paired regions were linked by an uneven number of unpaired bases.

Compared to TFs, which display complex dimerization patterns when bound to DNA, RBPs displayed simpler dimer spacing patterns. This is likely due to the fact that the backbone of a single-stranded nucleic acid has rotatable bonds. Thus, cooperativity between two RBDs requires that they bind to relatively closely spaced motifs.

Analysis of the biological roles of the RBP motif matches indicated that many motif matches were conserved, and specifically located at genomic features such as splice junctions. In particular, our analysis suggested a new role for ZC3H12, BOLL and DAZL proteins in regulating alternative splicing, and MEX3B in binding to type I interferon-regulated genes. As a large number of novel motifs were generated in the study, we expect that many other RBPs will have specific roles in particular biological functions.

332        Our results represent the largest single systematic study of human RNA-binding
333    proteins to date. This class of proteins is known to have major roles in RNA metabolism,
334    splicing and gene expression. However, the precise roles of RBPs in these biological
335    processes are poorly understood, and in general the field has been severely understudied.
336    The generated resource will greatly facilitate research in this important area.
337
338

339 **MATERIALS AND METHODS**
340
341
342 **Clone collection, cloning and protein production**
343
344 Clones were either collected from the human Orfeome 3.1 and 8.1 clone libraries (full
345 length clones) or ordered as synthetic genes from Genscript (eRBP constructs). As in our
346 previous work (Jolma et al., 2013), protein-coding synthetic genes or full length ORFs were
347 cloned into pETG20A-SBP to create an *E.coli* expression vector that allows the RBP or RBD
348 cDNAs to be fused N-terminally to Thioredoxin+6XHis and C-terminally to SBP-tags. Fusion
349 proteins were then expressed in the Rosetta P3 DE LysS *E.coli* strain (Novagen) using an
350 autoinduction protocol (Jolma et al., 2015). All constructs described in **Table S1** were
351 expressed and subjected to HTR-SELEX, regardless of protein level expressed. Protein
352 production was assessed in parallel by 96-well SDS-PAGE (ePage, Invitrogen). The success
353 rate of protein production was dependent on the size of the proteins, with most small RBDs
354 expressing well in *E.coli*. Significantly lower yield of protein was observed for full-length
355 proteins larger than 50 kDa.
356 After HIS-tag based IMAC purification, glycerol was added to a final concentration of
357 10%. Samples were split to single-use aliquots with approximately 200 ng RBP in a 5µl
358 volume and frozen at -80°C.
359
360
361 **Selection library generation**
362
363 To produce a library of RNA sequences for selection (selection ligands), we first
364 constructed dsDNA templates by combining three oligonucleotides together in a three cycle
365 PCR reaction (Phusion, NEB). For information about the barcoded ligand design, see **Table
366 S1**. The ligand design was similar to that used in our previous work analyzing TF binding
367 specificities in dsDNA (Jolma et al., 2013) except for the addition of a T7 RNA polymerase
368 promoter in the constant flanking regions of the ligand. RNA was expressed from the DNA-
369 templates using T7 *in vitro* transcription (Ampliscribe T7 High Yield Transcription Kit,
370 *Epicentre* or Megascript-kit *Ambion*) according to manufacturer's instructions, after which
371 the DNA-template was digested using RNAse-free DNAseI (Epicentre) or the TURBO-DNAse
372 supplied with the Megascript-kit. All RNA-production steps included RiboGuard RNAse-
373 inhibitor (Epicentre).
374 Two different approaches were used to facilitate the folding of RNA molecules. In the
375 protocol used in experiments where the batch identifier starts with letters "EM", RNA-
376 ligands were heated to +70°C followed by gradual, slow cooling to allow the RNA to fold into
377 minimal energy structures, whereas in batches "AAG" and "AAH" RNA transcription was not
378 followed by such folding protocol. The rationale was that spontaneous co-transcriptional
379 RNA-folding may better reflect folded RNA structures in the *in vivo* context. In almost all of
380 the cases where the same RBPs were tested with both of the protocols the results were highly
381 similar.
382
383

10

**HTR-SELEX assay**

Selection reactions were performed as follows: ~200ng of RBP was mixed on ice with ~1μg of the RNA selection ligands to yield approximate 1:5 molar ratio of protein to ligand in 20μl of Promega buffer (50 mM NaCl, 1 mM $MgCl_2$, 0.5 mM $Na_2EDTA$ and 4% glycerol in 50 mM Tris-Cl, pH 7.5). The complexity of the initial DNA library is approximately $10^{12}$ DNA molecules with 40 bp random sequence (~20 molecules of each 20 bp sequence on the top strand). The upper limit of detection of sequence features of HTR-SELEX is thus around 40 bits of information content.

The reaction was incubated for 15 minutes at +37°C followed by additional 15 minutes at room temperature in 96-well microwell plates (4-titude, USA), after which the reaction was combined with 50 μl of 1:50 diluted paramagnetic HIS-tag beads (His Mag Sepharose excel, GE-Healthcare) that had been blocked and equilibrated into the binding buffer supplemented with 0.1% Tween 20 and 0.1μg/μl of BSA (Molecular Biology Grade, NEB). Protein-RNA complexes were then incubated with the magnetic beads on a shaker for further two hours, after which the unbound ligands were separated from the bound beads through washing with a Biotek 405CW plate washer fitted with a magnetic platform. After the washes, the beads were suspended in heat elution buffer (0.5 μM RT-primer, 1 mM EDTA and 0.1% Tween20 in 10 mM Tris-Cl buffer, pH 7) and heated for 5 minutes at 70°C followed by cooling on ice to denature the proteins and anneal the reverse transcription primer to the recovered RNA library, followed by reverse transcription and PCR amplification of the ligands. The efficiency of the selection process was evaluated by running a qPCR reaction on parallel with the standard PCR reaction.


**Sequencing and generation of motifs**

PCR products from RNA libraries (indexed by bar-codes) were pooled together, purified using a PCR-purification kit (Qiagen) and sequenced using Illumina HiSeq 2000 (55 bp single reads). Data was de-multiplexed, and initial data analysis performed using the Autoseed algorithm (Nitta et al., 2015) that was further adapted to RNA analysis by taking into account only the transcribed strand and designating uracil rather than thymine. Autoseed identifies gapped and ungapped kmers that represent local maximal counts relative to similar sequences within their Huddinge neighborhood (Nitta et al., 2015). It then generates a draft motif using each such kmer as a seed. This initial set of motifs is then refined manually to identify the final seeds (**Table S2**), to remove artifacts due to selection bottlenecks and common "aptamer" motifs that are enriched by the HTR-SELEX process itself, and motifs that are very similar to each other. To assess initial data, we compared the deduced motifs to known motifs, to replicate experiments and experiments performed with paralogous proteins. Individual results that were not supported by replicate or prior experimental data were deemed inconclusive and were not included in the final dataset. Draft models were manually curated (AJ, JT, QM, TRH) to identify successful experiments, and final models were generated using the seeds indicated in **Table S2**.

Autoseed detected more than one seed for many RBPs. Up to four seeds were used to generate a maximum of two unstructured and two structured motifs. Of these, the motif with largest number of seed matches using the multinomial setting indicated on **Table S2** was

11

430 designated the primary motif. The motif with the second largest number of matches was
431 designated the secondary motif. The counts of the motifs represent the prevalence of the
432 corresponding motifs in the sequence pool (**Table S2**). Only these primary and secondary
433 motifs were included in further analyses. Such additional motifs are shown for LARP6 in **Fig.
434 S9**.
435   To find RBPs that bind to dimeric motifs, we visually examined the PWMs to find
436 direct repeat pattern of three or more base positions, with or without a gap between them
437 (see **Table S2**). The presence of such repetitive pattern could be either due to dimeric
438 binding, or the presence of two RBDs that bind to similar sequences in the same protein.
439   To identify structured motifs, we visually investigated the correlation diagrams for
440 each seed to find motifs that displayed the diagonal pattern evident in **Figure 2B**. The plots
441 display effect size and maximal sampling error, and show the deviation of nucleotide pair
442 distribution from what is expected from the distribution of the individual nucleotides. For
443 each structured motif, SLM models (**Table S3**) were built from sequences matching the
444 indicated seeds; a multinomial 2 setting was used to prevent the paired bases from
445 influencing each other. Specifically, when the number of occurrence of each pair of bases was
446 counted at the base-paired positions, neither of the paired bases was used to identify the
447 sequences that were analyzed. The SLMs were visualized either as the T-shaped logo (**Figure
448 3**) or as a PWM type logo where the bases that constitute the stem were shaded based on the
449 total fraction of A:U, G:C and G:U base pairs.
450   For analysis of RNA structure in **Figures 2** and **S5**, sequences matching the regular
451 expression NNNNCAGU[17N]AGGCNNN or sequences of the three human collagen gene
452 transcripts (From 5' untranslated and the beginning  the coding sequence, the start codon is
453 marked with bold typeface: COL1A1 -CCACAAAGAGUCUACAUGUCUAGGGUCUAG-
454 AC**AUG**UUCAGCUUUGUGG; COL1A2- CACAAGGAGUCUGCAUGUCUAAGUGCUAGA-
455 C**AUG**CUCAGCUUUGUG and COL3A1 - CCACAAAGAGUCUACAUGGGUC**AUG**UUCAG-
456 CUUUGUGG) were analysed using "RNAstructure" software (Mathews, 2014) through the
457 web-interface
458 in:http://rna.urmc.rochester.edu/RNAstructureWeb/Servers/Fold/Fold.html using default
459 settings. All structures are based on the program's minimum energy structure prediction.
460 For analysis in **Figure 3**, we extracted all sequences that matched the binding sequences of
461 MKRN1 and ZRANB2 (GUAAAKUGUAG and NNNGGUAAGGUNN, respectively; N denotes a
462 weakly specified base) flanked with ten bases on both sides from the cycle four of HTR-
463 SELEX. Subsequently, we predicted their secondary structures using the program RNAfold
464 (Vienna RNA package; (Lorenz et al., 2011)) followed by counting the predicted secondary
465 structure at each base position in the best reported model for each sequence. For both RBPs,
466 the most common secondary structure for the bases within the defined part of the consensus
467 (GUAAAKUGUAG and GGUAAGGU) was the fully single stranded state (82% and 30% of all
468 predicted structures, respectively). To estimate the secondary structure at the flanks, the
469 number of paired bases formed between the two flanks were identified for each sequence.
470 Fraction of sequences with specific number of paired bases are shown in **Figure 3**.
471
472
473
474
475

12

**Motif mapping**

To gain insight into the function of the RBPs, we mapped each motif to the whole human genome (hg38). We applied different strategies for the linear and the stem-loop motifs. For the linear motifs, we identified the motif matches with MOODS (Korhonen et al., 2017) with the following parameter setting: --best-hits 300000 --no-snps. For the stem-loop motifs, we implemented a novel method to score sequences against the SLMs. The source code is available on GitHub: https://github.com/zhjilin/rmap.

We identified the 300,000 best scored matches in the genome, and further included any matches that had the same score as the match with the lowest score, leading to at least 300,000 matches for each motif. The matches were then intersected with the annotated features from the ENSEMBL database (hg38, version 91), including the splicing donor (DONOR), splicing acceptor (ACCEPTOR), the translation start codon (STARTcodon), the translation stop codon (STOPcodon) and the transcription starting site (TSS). The above features were filtered in order to remove short introns (<50bp), and features with non-intact or non-canonical start codon or stop codon. The filtered features were further extended 1kb both upstream and downstream in order to place the feature in the centre of all the intervals. The motif matches overlapping the features were counted using BEDTOOLS (version 2.15.0) and normalized by the total number of genomic matches for the corresponding motif.


**Motif comparisons and GO analysis**

To assess the similarity between publicly available motifs and our HTR-SELEX data, we aligned the motifs as described in (Jolma et al., 2015) (**Figure S1**). To determine whether RBPs with similar RBDs recognize and bind to similar targets, we compared the sequences of the RBDs and their motifs. First, the RBPs were classified based on the type and number of RBDs. For each class, we then extracted the amino-acid sequence of the RBPs starting from the first amino acid of the first RBD and ending at the last amino acid of the last RBD. We also confirmed the annotation of the RBDs by querying each amino acid sequence against that SMART database, and annotated the exact coordinates of the domains through the web-tools: http://smart.embl-heidelberg.de and http://smart.embl-heidelberg.de/smart/batch.pl. Sequence similarities and trees were built using PRANK (Loytynoja and Goldman, 2005) (parameters: -d, -o, -showtree). The structure of the tree representing the similarity of the domain sequence was visualized using R (version 3.3.1).

For identification of classes of transcripts that are enriched in motif matches for each RBP, we extracted the top 100 transcripts according to the score density of each RBP motif. These 100 transcripts were compared to the whole transcriptome to conduct the GO enrichment analysis for each motif using the R package ClusterProfiler (version 3.0.5).

To analyze conservation of motif matches, sites recognized by each motif were searched from both strands of 100 bp windows centered at the features of interest (acceptor, donor sites) using the MOODS program (version 1.0.2.1). For each motif and feature type, 1000 highest affinity sites were selected for further analysis regardless of the matching strand. Whether the evolutionary conservation of the high affinity sites was explained by the motifs was tested using program SiPhy (version 0.5, task 16, seedMinScore 0) and multiz100way multiple alignments of 99 vertebrate species to human (downloaded from

13

522    UCSC genome browser, version hg38). A site was marked as being conserved according to
523    the motif if its SiPhy score was positive meaning that the aligned bases at the site were better
524    explained by the motif than by a neutral evolutionary model (hg38.phastCons100way.mod
525    obtained from UCSC genome browser). Two motifs were excluded from the analysis because
526    the number of high affinity sites that could be evaluated by SiPhy was too small. The
527    hypothesis that the motif sites in the sense strand were more likely to be conserved than
528    sites in the antisense strand was tested against the null hypothesis that there was no
529    association between site strand and conservation using Fisher's exact test (one-sided). The
530    P values given by the tests for individual motifs were corrected for multiple testing using
531    Holm's method.
532
533
534    **Mutual information calculation**
535
536         The global pattern of motifs across the features tested was analyzed by calculating
537    the mutual information (MI) between 3-mer distributions at two non-overlapping positions
538    of the aligned RNA sequences. MI can be used for such analysis, because if a binding event
539    contacts two continuous or spaced 3-bp wide positions of the sequences at the same time,
540    the 3-mer distributions at these two positions will be correlated. Such biased joint
541    distribution can then be detected as an increase in MI between the positions.
542         Specifically, MI between two non-overlapping positions (pos1, pos2) was estimated
543    using the observed frequencies of a 3-mer pair (3+3-mer), and of its constituent 3-mers at
544    both positions:

$$MI(pos1, pos2) = \sum P(3\text{+}3\text{-}mer)\log_2 \frac{P(3\text{+}3\text{-}mer)}{P_{pos1}(3\text{-}mer)P_{pos2}(3\text{-}mer)}$$

545    where $P(3\text{+}3\text{-}mer)$ is the observed probability of the 3-mer pair (i.e. gapped or ungapped 6
546    mer). $P_{pos1}(3\text{-}mer)$ and $P_{pos2}(3\text{-}mer)$, respectively, are the marginal probabilities of the
547    constitutive 3-mers at position 1 and position 2. The sum is over all possible 3-mer pairs.
548         To focus on RBPs that specifically bind to a few closely related sequences, such as
549    RBPs with well-defined motifs, it is possible to filter out most background non-specific
550    bindings (e.g., selection on the shape of RNA backbone) by restricting the MI calculation, to
551    consider only the most enriched 3-mer pairs for each two non-overlapping positions.
552         Such enriched 3-mer pair based mutual information (E-MI) is calculated by summing
553    MI over top-10 most enriched 3-mer pairs.
554

$$E\text{-}MI(pos1, pos2) = \sum_{top\ 3\text{+}3\text{-}mers} P(3\text{+}3\text{-}mer)\log_2 \frac{P(3\text{+}3\text{-}mer)}{P_{pos1}(3\text{-}mer)P_{pos2}(3\text{-}mer)}$$

555
556
557
558
559

14

560 **FIGURE LEGENDS**
561
562 **Figure 1. RNA HT-SELEX protocol and data-analysis**
563 (**A**) Schematic illustration of the HTR-SELEX process. RBD or full-length RBPs expressed in
564 *E.coli* as HIS$_6$-tagged fusion proteins (left top) were purified and incubated with barcoded
565 RNA selection ligands. RNA ligands bound by the proteins were recovered by RT-PCR,
566 followed by *in vitro* transcription to generate the RNA for the next cycle of SELEX (left
567 middle). The procedure was repeated at least three times and the ligands recovered from the
568 selection cycles were subjected to Illumina sequencing (left bottom) with data analysis to
569 generate binding specificity models (right).
570 (**B**) Comparison of the number of RBPs with motifs derived in the present study (HTR-SELEX)
571 and two previous studies (RNAcompete (Cook et al., 2011) and SELEX (Ray et al., 2013)).
572 Note that our analysis revealed motifs for 49 RBPs for which a motif was not known.
573 (**C**) Distribution of RBP with motifs classified by structural family of the RBDs.
574
575 **Figure 2. Detection of RNA binding models**
576 (**A**) ZRANB2 binds to a linear RNA motif. The motif of ZRANB2 is shown below the triangular
577 correlation heatmap. The heatmap illustrates deviation of the observed nucleotide
578 distributions from those predicted by a mononucleotide model where bases are independent.
579 (**B**) MKRN1 binds preferentially to a stem-loop. Note a diagonal series of red tiles (boxed)
580 that indicates pairs of bases whose distribution deviates from the independence assumption.
581 These bases are shaded in the motif below the triangle. The interdependency occurs between
582 bases that are at the same distance from the center of the motif, consistent with formation of
583 a stem-loop structure. Right top: A RNAfold-predicted stem-loop structure for a sequence
584 that was highly enriched in the experiment.
585 (**C**) LARP6 binds to a complex, bulged RNA structure. The left panel indicates the
586 dinucleotide dependencies, whereas the right panel presents a predicted structure of the
587 bound RNA.
588 (**D**) Fraction of RBPs with linear and structured binding specificities.
589 (**E**) Length distribution of stem and loop for the structured motifs.
590
591 **Figure 3. Comparison between linear PWM and stem loop (SLM) models.**
592 Left: Visualization of the stem loop models. A T-shape model (top) shows a horizontal loop
593 and a vertical stem where the frequency of each base combination is shown. Bases are
594 aligned so that Watson-Crick base pairs orient horizontally. Pie-charts show frequency of
595 Watson-Crick (green) and G-U (light green) at each position of the predicted stem. A linear
596 visualization (bottom) where the base pairing frequency is indicated by the darkness of gray
597 shading is also shown.
598 Middle: Schematic description of the scoring process for the SLM. All possible alignment
599 positions between an 8-mer with a 4 base gap in the middle and the model are searched in
600 order to find the aligned position with the best score. When the 8-mer overlaps both bases
601 of a SLM-predicted base-pair, the score for the paired position (red tiles connected by black
602 lines) is derived from the SLM base-pair score. In cases where the kmer aligns to only one
603 base of the SLM base-pair, the score for the position (black) is derived from the
604 mononucleotide matrix.

15

Right: RNA secondary structure prediction analysis using RNAfold reveals that sequences flanking MKRN1 loop sequence form base pairs (top), whereas bases on the flanks of ZRANB2 matches (bottom) are mostly unpaired.

**Figure 4. Comparison between the obtained motifs**
(**A** to **H**) Similar RBPs bind to similar motifs. Motifs were classified into six major categories based on structural class of the RBPs. Dendrograms are based on amino-acid alignment using PRANK. Within RRM family, RBPs with a distinct number of RRMs were grouped and aligned separately. Motifs shown are the primary motif for each RBP. Asterisks indicate a stem-loop structured motif, with the gray shading showing the strength of the base pairing at the corresponding position. Two asterisks indicate that the RBP can binds to a structured secondary motif. Only families with more than one representative RHT-SELEX motif are shown.
(**G**) RBPs commonly prefer sequences with G or U nucleotides. Frequencies of all mononucleotides (left) and dinucleotides (right) across all of the RBP motifs. Note that G and U are overrepresented.

**Figure 5. Many RBPs can recognize more than one motif**
Pie chart (top) indicates fraction of RBPs that recognize more than one motif. Primary (left) and secondary (right) motifs are shown, classified according to the RBP structural family. Right column indicates whether the RBP binds to two linear motifs (+), two structured motifs (red circle) or both types of motifs (green circle).

**Figure 6. RBP motif matches are conserved and enriched in distinct sequence features**
(**A**) Strong enrichment of RBP motif matches at or near the splicing donor and acceptor sites. Mononucleotide frequencies at splice donor and acceptor sites are shown on top, above the gene schematic. Left: meta-plots indicate the enrichment of ZRANB2 and ZC3H12C motif matches at splice donor sites. Right: enrichment of BOLL and DAZL at splice acceptor sites. Blue dots indicate the number of matches in the sense strand at each base position; black line indicates the median in 10 base sliding windows. Corresponding values for the anti-sense strand are shown as light blue dots and dotted black line, respectively.
(**B**) The conservation of motif binding sites in sense vs. antisense strand. For each feature type (acceptor, donor site) and binding motif, thousand highest affinity sites in total were selected from the hundred base windows centered at the features. The total number the conserved sites (x-axis) is plotted against odds ratio of conserved vs. non-conserved site being in sense strand (y-axis). Those motifs for which conservation was significantly associated with sense strand (one-sided Fisher's exact test) are shown in green. The five motifs with smallest p-values are named.
(**C**) Enrichment of ZRANB2 and ZC3H12C motif matches near transcription start sites (TSS). Note that matches are only enriched on the sense strand downstream of the TSS.
(**D**) Gene Ontology enrichment of MEX3B motif matches. The top 100 genes with highest motif-matching score density were used to conduct the Gene Ontology enrichment analysis. The enriched GO terms were simplified by their similarity (cutoff=0.5). The fraction of genes in the GO categories is also shown (Gene Ratio).

16

651    **Supplementary information**
652
653    **Figure S1. The similarity of motifs between HTR-SELEX and the curated dataset**. In total,
654    33 motifs from the CISBP-RNA database (RNAcompete) were collected to compare with the
655    HTR-SELEX derived motifs (both primary and secondary motifs). Motifs were presented and
656    organized according to their protein structure family by descending similarity score. Higher
657    score indicates higher similarity between motifs.
658
659    **Figure S2. Information content and width distribution of HTR-SELEX motifs.** The per-
660    base information was calculated for every individual position in the PWM. The overall
661    information content of each motif is the sum of all positions in the PWM. The width of each
662    motif was generated by counting the number of position in the corresponding PWM.
663
664    **Figure S3. RBPs with multimeric binding sites.**  Some RBPs (31.3%, left) bind to the
665    sequence as homodimers. Two identical half-sites are separated by a spacing sequence. The
666    distribution of spacing preference of all RBPs is shown (right).
667
668    **Figure S4. Spacing preferences between dimeric binding sites are consistent in**
669    **different assays.** For four RBPs, the same seeds were used in different assays to detect the
670    spacing preferences. The colour-coded arrays represent the spacing information extracted
671    from HTR-SELEX, PAR-CLIP and HITS-CLIP. The results are consistent between HTR-SELEX
672    (top row) and PAR-CLIP (bottom row). Pie charts show the percentage of reads containing
673    the indicated instances of the motifs CAC in RBPMS binding motifs as determined by HTR-
674    SELEX (left) or PAR-CLIP (middle). Incidence of the motif CCA is also shown (right).
675
676    **Figure S5. Known binding motifs of LARP6**. The left three structures were generated using
677    the sequences enriched in HTR-SELEX. The right three structures illustrate the predicted
678    structures of known collagen RNA sequences.
679
680    **Figure S6. RBP families with and without structural specificity.** The count of RBPs
681    recognizing structured and unstructured binding motifs in each protein structure family.
682
683    **Figure S7. Information content correlation between the SLM and the mono-nucleotide**
684    **PWM. Left. Information content correlation per base. Right. Overall information**
685    **content correlation.** In general, the SLM yielded higher per-base information content.
686
687    **Figure S8. Dominating set of HTR-SELEX motifs.** Cystoscope (Version 3.2.1) was used to
688    visualize the dominating set on top of the relationship map between motifs with a cutoff of
689    5e-6, calculated by SSTAT (see the method part).
690
691    **Figure S9. Various binding specificities detected in LARP6.** LARP6 is able to recognise
692    and bind to distinct sequences through different strategies besides binding to the internal
693    loop structure.  (**A**) Short and long linear motifs (**B**) unstructured motifs with gaps.
694

695   **Figure S10. The mutual information (MI) meta-plots around the splicing donor and**
696   **acceptor sites.** The splice donor and acceptor sites are placed in the centre of the 147nts
697   sequence. The detected signals close to the donor and acceptor sites are shown in red.
698
699   **Table S1. Sequence information of proteins and DNA.**
700   **Table S2. PWMs of the linear motifs.**
701   **Table S3. PWMs of the structured motifs**
702   **Table S4. Dependency matrices of paired bases for the structured motifs.**
703
704   **Supplementary Data S1. Meta-plots of the motif matches enrichment at splice donor,**
705   **acceptor, TSS, start and stop codon.**
706

707 **Reference**

708

709 Bakheet, T., Hitti, E., Al-Saif, M., Moghrabi, W.N., and Khabar, K.S.A. (2018). The AU-rich
710 element landscape across human transcriptome reveals a large proportion in introns and
711 regulation by ELAVL1/HuR. Biochim Biophys Acta *1861*, 167-177.
712 Cai, L., Fritz, D., Stefanovic, L., and Stefanovic, B. (2010). Binding of LARP6 to the conserved
713 5' stem-loop regulates translation of mRNAs encoding type I collagen. J Mol Biol *395*, 309-
714 326.
715 Chen, Y., Zubovic, L., Yang, F., Godin, K., Pavelitz, T., Castellanos, J., Macchi, P., and Varani, G.
716 (2016). Rbfox proteins regulate microRNA biogenesis by sequence-specific binding to their
717 precursors and target downstream Dicer. Nucleic Acids Res *44*, 4381-4395.
718 Cook, K.B., Kazan, H., Zuberi, K., Morris, Q., and Hughes, T.R. (2011). RBPDB: a database of
719 RNA-binding specificities. Nucleic Acids Research *39*, D301-D308.
720 Damianov, A., Kann, M., Lane, W.S., and Bindereif, A. (2006). Human RBM28 protein is a
721 specific nucleolar component of the spliceosomal snRNPs. Biol Chem *387*, 1455-1460.
722 Darnell, R.B. (2010). HITS-CLIP: panoramic views of protein-RNA regulation in living cells.
723 Wiley interdisciplinary reviews RNA *1*, 266-286.
724 Farazi, T.A., Leonhardt, C.S., Mukherjee, N., Mihailovic, A., Li, S., Max, K.E., Meyer, C., Yamaji,
725 M., Cekan, P., Jacobs, N.C.*, et al.* (2014). Identification of the RNA recognition element of the
726 RBPMS family of RNA-binding proteins and their transcriptome-wide mRNA targets. RNA
727 (New York, NY) *20*, 1090-1102.
728 Gerstberger, S., Hafner, M., and Tuschl, T. (2014). A census of human RNA-binding proteins.
729 Nature reviews Genetics *15*, 829-845.
730 Glisovic, T., Bachorik, J.L., Yong, J., and Dreyfuss, G. (2008). RNA-binding proteins and post-
731 transcriptional gene regulation. FEBS letters *582*, 1977-1986.
732 Hafner, M., Landthaler, M., Burger, L., Khorshid, M., Hausser, J., Berninger, P., Rothballer, A.,
733 Ascano, M., Jr., Jungkamp, A.C., Munschauer, M.*, et al.* (2010). Transcriptome-wide
734 identification of RNA-binding protein and microRNA target sites by PAR-CLIP. Cell *141*,
735 129-141.
736 Haque, N., Ouda, R., Chen, C., Ozato, K., and Hogg, J.R. (2018). ZFR coordinates crosstalk
737 between RNA decay and transcription in innate immunity. Nat Commun *9*, 1145.
738 Hayakawa-Yano, Y., Suyama, S., Nogami, M., Yugami, M., Koya, I., Furukawa, T., Zhou, L., Abe,
739 M., Sakimura, K., Takebayashi, H.*, et al.* (2017). An RNA-binding protein, Qki5, regulates
740 embryonic neural stem cells through pre-mRNA processing in cell adhesion signaling.
741 Genes Dev *31*, 1910-1925.
742 Huang, H., Weng, H., Sun, W., Qin, X., Shi, H., Wu, H., Zhao, B.S., Mesquita, A., Liu, C., Yuan,
743 C.L.*, et al.* (2018). Recognition of RNA N(6)-methyladenosine by IGF2BP proteins enhances
744 mRNA stability and translation. Nat Cell Biol *20*, 285-295.
745 Jolma, A., Yan, J., Whitington, T., Toivonen, J., Nitta, K.R., Rastas, P., Morgunova, E., Enge, M.,
746 Taipale, M., Wei, G.*, et al.* (2013). DNA-binding specificities of human transcription factors.
747 Cell *152*, 327-339.
748 Jolma, A., Yin, Y., Nitta, K.R., Dave, K., Popov, A., Taipale, M., Enge, M., Kivioja, T., Morgunova,
749 E., and Taipale, J. (2015). DNA-dependent formation of transcription factor pairs alters
750 their binding specificity. Nature *527*, 384-388.
751 Korhonen, J.H., Palin, K., Taipale, J., and Ukkonen, E. (2017). Fast motif matching revisited:
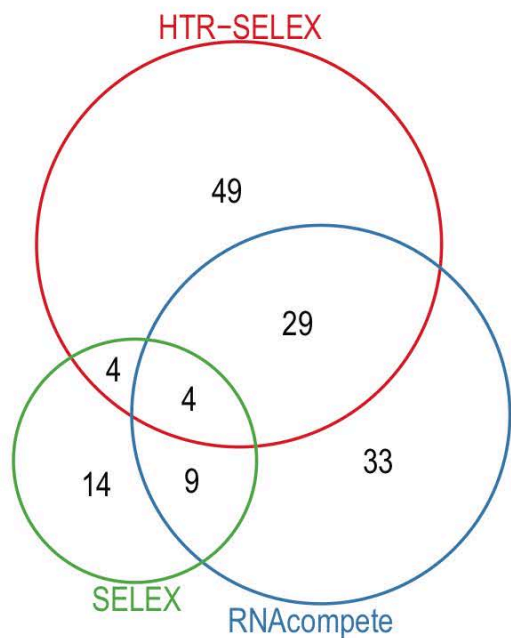752 high-order PWMs, SNPs and indels. Bioinformatics *33*, 514-521.

Lambert, N.J., Robertson, A.D., and Burge, C.B. (2015). RNA Bind-n-Seq: Measuring the Binding Affinity Landscape of RNA-Binding Proteins. Methods in enzymology *558*, 465-493.

Lamesch, P., Li, N., Milstein, S., Fan, C., Hao, T., Szabo, G., Hu, Z., Venkatesan, K., Bethel, G., Martin, P.*, et al.* (2007). hORFeome v3.1: a resource of human open reading frames representing over 10,000 human genes. Genomics *89*, 307-315.

Leppek, K., Schott, J., Reitter, S., Poetz, F., Hammond, M.C., and Stoecklin, G. (2013). Roquin promotes constitutive mRNA decay via a conserved class of stem-loop recognition motifs. Cell *153*, 869-881.

Lorenz, R., Bernhart, S.H., Honer Zu Siederdissen, C., Tafer, H., Flamm, C., Stadler, P.F., and Hofacker, I.L. (2011). ViennaRNA Package 2.0. Algorithms Mol Biol *6*, 26.

Loughlin, F.E., Mansfield, R.E., Vaz, P.M., McGrath, A.P., Setiyaputra, S., Gamsjaeger, R., Chen, E.S., Morris, B.J., Guss, J.M., and Mackay, J.P. (2009). The zinc fingers of the SR-like protein ZRANB2 are single-stranded RNA-binding domains that recognize 5' splice site-like sequences. Proc Natl Acad Sci U S A *106*, 5581-5586.

Loytynoja, A., and Goldman, N. (2005). An algorithm for progressive multiple alignment of sequences with insertions. Proc Natl Acad Sci U S A *102*, 10557-10562.

Mathews, D.H. (2014). RNA Secondary Structure Analysis Using RNAstructure. Curr Protoc Bioinformatics *46*, 12 16 11-25.

Nitta, K.R., Jolma, A., Yin, Y., Morgunova, E., Kivioja, T., Akhtar, J., Hens, K., Toivonen, J., Deplancke, B., Furlong, E.E.*, et al.* (2015). Conservation of transcription factor binding specificities across 600 million years of bilateria evolution. Elife *4*.

Orenstein, Y., Ohler, U., and Berger, B. (2018). Finding RNA structure in the unstructured RBPome. BMC Genomics *19*, 154.

Ray, D., Kazan, H., Chan, E.T., Pena Castillo, L., Chaudhry, S., Talukder, S., Blencowe, B.J., Morris, Q., and Hughes, T.R. (2009). Rapid and systematic analysis of the RNA recognition specificities of RNA-binding proteins. Nat Biotechnol *27*, 667-670.

Ray, D., Kazan, H., Cook, K.B., Weirauch, M.T., Najafabadi, H.S., Li, X., Gueroussov, S., Albu, M., Zheng, H., Yang, A.*, et al.* (2013). A compendium of RNA-binding motifs for decoding gene regulation. Nature *499*, 172-177.

Rosario, R., Childs, A.J., and Anderson, R.A. (2017). RNA-binding proteins in human oogenesis: Balancing differentiation and self-renewal in the female fetal germline. Stem Cell Res *21*, 193-201.

Skrisovska, L., Bourgeois, C.F., Stefl, R., Grellscheid, S.N., Kister, L., Wenter, P., Elliott, D.J., Stevenin, J., and Allain, F.H. (2007). The testis-specific human protein RBMY recognizes RNA through a novel mode of interaction. EMBO Rep *8*, 372-379.

Tuerk, C., and Gold, L. (1990). Systematic evolution of ligands by exponential enrichment: RNA ligands to bacteriophage T4 DNA polymerase. Science (New York, NY) *249*, 505-510.

Varani, G., and McClain, W.H. (2000). The G x U wobble base pair. A fundamental building block of RNA structure crucial to RNA function in diverse biological systems. EMBO Rep *1*, 18-23.

Weyn-Vanhentenryck, S.M., Mele, A., Yan, Q., Sun, S., Farny, N., Zhang, Z., Xue, C., Herre, M., Silver, P.A., Zhang, M.Q.*, et al.* (2014). HITS-CLIP and integrative modeling define the Rbfox splicing-regulatory network linked to brain development and autism. Cell reports *6*, 1139-1152.

797   Xia, H., Chen, D., Wu, Q., Wu, G., Zhou, Y., Zhang, Y., and Zhang, L. (2017). CELF1
798   preferentially binds to exon-intron boundary and regulates alternative splicing in HeLa
799   cells. Biochim Biophys Acta *1860*, 911-921.
800   Yin, Y., Morgunova, E., Jolma, A., Kaasinen, E., Sahu, B., Khund-Sayeed, S., Das, P.K., Kivioja,
801   T., Dave, K., Zhong, F.*, et al.* (2017). Impact of cytosine methylation on DNA binding
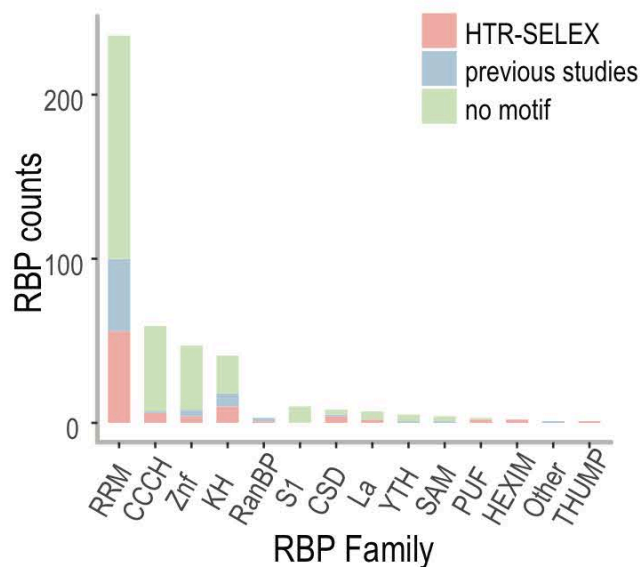802   specificities of human transcription factors. Science (New York, NY) *356*.
803

# Jolma et al. 2018 Figure 1



**A** *HT-SELEX process*

~950 RBP clones    Synthetic ligands

*E.coli*    RBPs    RNAs

*Library preparation*

Sequencing library

*Sequencing*

Reads

*Motif discovery*

HEXIM1
CARHSP1
CELF1
ZC3H12B
THUMPD1
QKI

100
50
0
Watson-Crick & G-U base-pairing (%)

**B** *Comparison to other datasets*

HTR−SELEX

49

29

4

4

14    9    33

SELEX    RNAcompete

**C** *RBPs based on protein domain*

HTR-SELEX
previous studies
no motif

RBP counts

200

100

0

RRM CCCH Znf KH RanBP S1 CSD La YTH SAM PUF HEXIM Other THUMP

RBP Family

# Jolma et al. 2018 Figure 2

**A** *Detection of linear motifs*



**B** *Detection of structured motifs*



**C** *Recognition of bulged stem-loop structures*



**D** *RBPs with structured motifs*



**E** *Length distribution of stem and loop*

# Jolma et al. 2018 Figure 3

*T shape illustration of stem-loop model*



Watson-Crick
UG + GU
others

position of scoring matrix

0 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15

Center

gapped kmer

*Motif with stem shaded*

base scored with mononucleotide matrix
base scored with SLM base-pair score
base without overlap with the score matrix

gap in the kmer
indicates base-pair

**MKRN1**

5 Paired bases 7%
4 Paired bases 4%
3 Paired bases 1%
2 or less paired bases 0%
6 Paired bases 15%
7 Paired bases 23%
8 Paired bases 26%
9 Paired bases 20%
10 Paired bases 4%

**ZRANB2**

10 Paired bases 0%
9 Paired bases 1%
8 Paired bases 2%
7 Paired bases 3%
6 Paired bases 4%
5 Paired bases 3%
4 Paired bases 2%
3 Paired bases 0%
2 or less paired bases 85%

# Jolma et al. 2018 Figure 4

# Jolma et al. 2018 Figure 5

# Jolma et al. 2018 Figure 6

# Jolma et al. 2018 Figure S1



| | | HTR–SELEX | | RNAcompete | Similarity score | |
|---|---|---|---|---|---|---|
| | | Primary | Secondary | | Primary | Secondary |
| **RRM** | RBFOX1 | | | | 10.17 | |
| | RBM42 | | | | 9.4 | |
| | SNRPA | | | | 7.98 | 7.91 |
| | RBM4 | | | | 7.96 | 3.17 |
| | PABPC5 | | | | 7.92 | 4.71 |
| | SNRNP70 | | | | 7.8 | 7.07 |
| | ZCRB1 | | | | 7.68 | 7.66 |
| | HNRNPA1L2 | | | | 7.17 | |
| | HNRNPCL1 | | | | 6.85 | |
| | HNRNPA1 | | | | 6.7 | |
| | DAZAP1 | | | | 5.39 | |
| | HNRNPC | | | | 5.1 | |
| | RBM24 | | | | 4.88 | |
| | CELF4 | | | | 4.41 | |
| | PCBP1 | | | | 4.26 | |
| | ELAVL1 | | | | 3.67 | |
| | MSI1 | | | | 3.45 | |
| | HNRNPL | | | | 2.77 | |
| | RBMS3 | | | | 1.38 | 2.17 |
| | RBM46 | | | | 1.1 | |
| | RBMS1 | | | | 0.96 | 3.96 |
| | RBM28 | | | | −1 | 0.21 |
| | RALY | | | | −1.24 | |
| | RBM6 | | | | −1.54 | −2.71 |
| | SART3 | | | | −3.24 | |
| **ZNF** | ZC3H10 | | | | 3.42 | 5.41 |
| | TARDBP | | | | 3.88 | 3.66 |
| **KH** | QKI | | | | 9.86 | |
| | KHDRBS1 | | | | 6 | |
| | KHDRBS2 | | | | 4.24 | |
| | KHDRBS3 | | | | 3.25 | |
| **CSD** | YBX1 | | | | 5.06 | 6.16 |
| | YBX2 | | | | 7.26 | 3.01 |

# Jolma et al. 2018 Figure S2

# Jolma et al. 2018 Figure S3



RBPs with dimeric binding specificities

Spacing of dimeric motifs

# Jolma et al. 2018 Figure S4



**PUM2 UGUA-N-UGUA**
HTR-SELEX
PAR-CLIP

Hafner et al. 2010

**IGF2BP1 CA[GU]-N-CA[GU]**
HTR-SELEX
PAR-CLIP

Hafner et al. 2010

**RBFOX1 GCAUG-N-GCAUG**
HITS-CLIP
HTR-SELEX
PAR-CLIP

Vanhentenryck et al. 2014

**RBPMS CAC-N-CAC**
HTR-SELEX
PAR-CLIP

Farazi et al. 2014

**Instances of CAC in RBPMS HTR-SELEX reads**

>5 1%
4 8%
0 11%
1 14%
2 37%
3 29%

**Instances of CCA in RBPMS HTR-SELEX reads**

>2 5%
1 27%
0 68%

**Instances of CAC in RBPMS PAR-CLIP clusters**

>5 1%
0 3%
4 7%
1 24%
2 40%
3 25%

# Jolma et al. 2018 Figure S5



Probability >= 99%
99% > Probability >= 95%
95% > Probability >= 90%
90% > Probability >= 80%
80% > Probability >= 70%
70% > Probability >= 60%
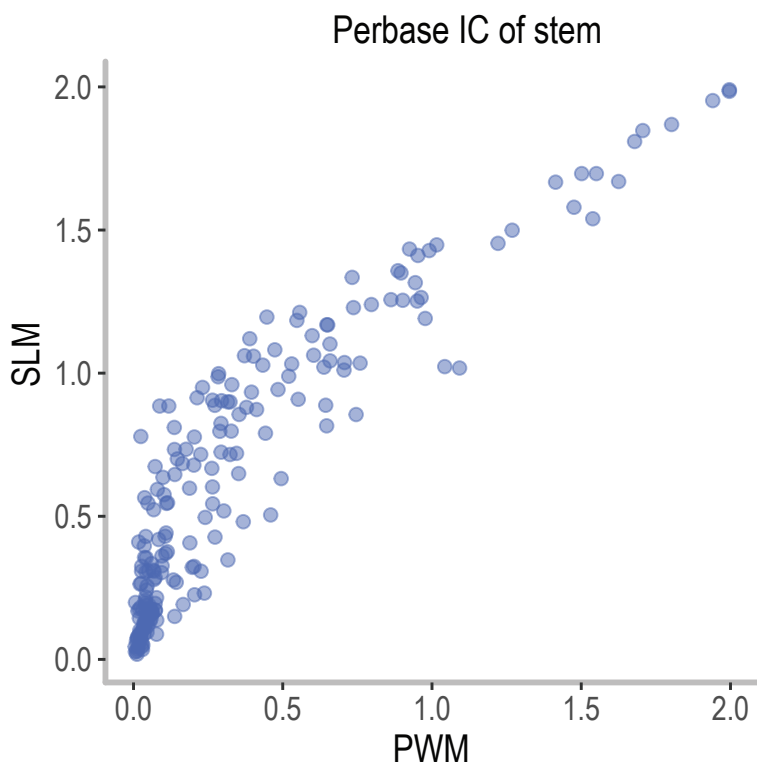60% > Probability >= 50%
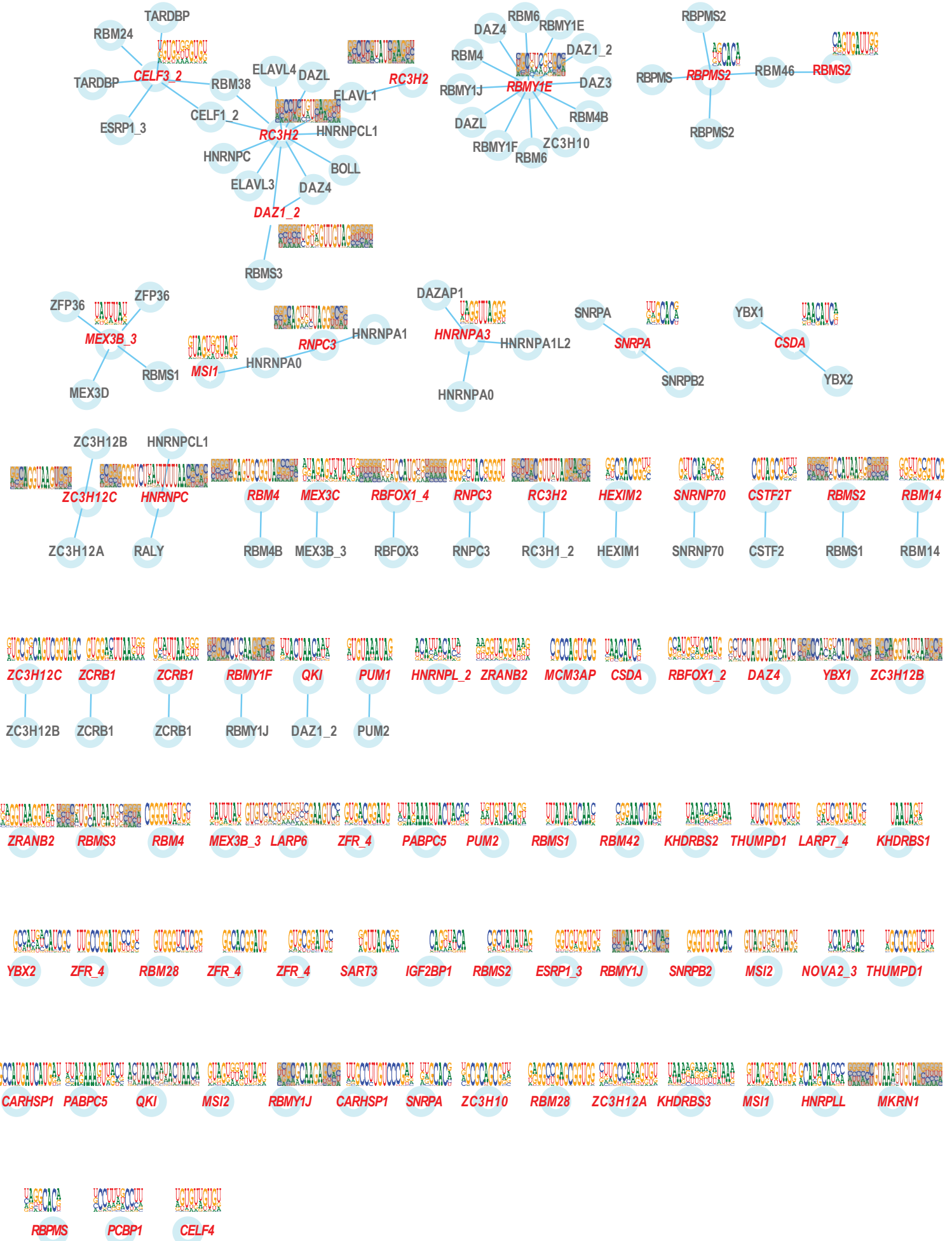50% > Probability

# Jolma et al. 2018 Figure S6

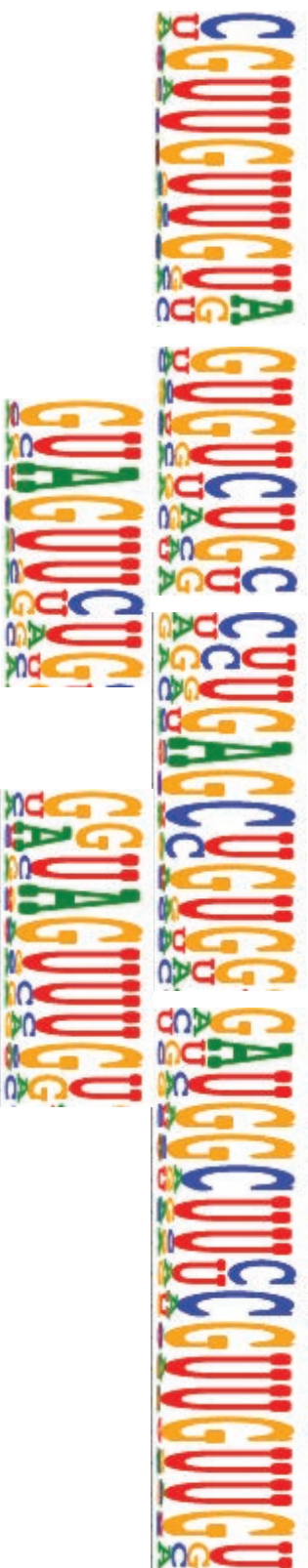# Jolma et al. 2018 Figure S7



Perbase IC of stem

Overall IC of stem

# Jolma et al. 2018 Figure S8

# Jolma et al. 2018 Figure S9

**A** *Simple non-structured motifs*

**B** *Variably gapped non-structured motifs*

GUGUCNGAAGU

UGGGCNGAGCC

GUGUCNCUAGC

# Jolma et al. 2018 Figure S10