

# **Evolutionary history of Cytoplasmic Polyadenylation Element-Binding Proteins**

**Principal Investigator: Labib Rouhana and Joseph Ryan**

**Draft or Version Number: v.1.1**

**19 December 2020**

## LIST OF ABBREVIATIONS

|               |   |
|---------------|---|
| <Insert text> | <Insert text>                                       |
| CPEB          | Cytoplasmic Polyadenylation Element-Binding Protein |
| RBP           | RNA-Binding Protein                                 |
| RRM           | RNA Recognition Motif                               |
| RBD           | RNA-Binding Domain (region of CPEB containing RRM)  |
| ZZ            | Zinc (binding domain)                               |
| poly(A)       | Polyadenosine                                       |

## 1 INTRODUCTION: BACKGROUND INFORMATION AND SCIENTIFIC RATIONALE

## 1.1 BACKGROUND INFORMATION

Cytoplasmic Polyadenylation Element-Binding Proteins (CPEBs) are pivotal regulators of gene expression during gametogenesis and the initial stages of embryonic development of bilaterians. These RNA-binding proteins recognize sequences (consensus UUUUA<sub>1-2</sub>U) in the 3'untranslated region of specific mRNAs and recruit cytoplasmic poly(A)-polymerases that increase the poly(A)-tail length of substrate mRNAs and consequently their translation. CPEBs have not been studied in outside of bilaterians, but a recent study identified them as one of "25 groups of metazoan-specific genes that are essential across the Animal Kingdom" (Paps and Holland, 2018).

## 1.2 RATIONALE

CPEB orthologs corresponding to both CPEB1 and CPEB2 subgroups have been found in many bilaterian groups, but detailed phylogenetic analyses that have included non-bilaterian animals as well as closely related non-animal lineages have yet to be performed. The results of such analyses could shine light on the origin and evolution of cytoplasmic polyadenylation.

## 1.3 HYPOTHESES

Orthologs of CPEB1 and CPEB2 are present in the genomes of all animals and absent from the genomes of non-animals. The most identical proteins in non-animal genomes will be RNA-binding proteins with essential functions shared throughout Eukaryota. Metazoan proteins with sequence similarity with CPEBs, which are not CPEB orthologs or orthologs of proteins present in non-metazoan genomes, share biochemical and biological functions with CPEBs.

## 1.4 OBJECTIVES

1. To validate the presence or absence of CPEB orthologs in genomes of non-bilaterian animals.
2. To identify non-animal proteins that share a most-recent common genic ancestor with CPEBs.
3. To document the evolutionary history of CPEBs and related genes.

## 2 STUDY DESIGN AND ENDPOINTS

- 1) **Identify putative CPEBs within genomes:** Starting with the RBD (tandem RRM and ZZ domain) of human CPEB1 (NP\_001275748.1; AA 234-479) and CPEB4 (AAH36899.1; AA 66-309) protein sequences (Afroz et al., 2014), perform TBLASTN searches against *Saccharomyces cerevisiae* and *Schizosaccharomyces pombe* (Fungi), *Arabidopsis thaliana*, *Capsaspora oewazarki* (Filasterea), *Salpingoeca*

*rosetta* (Choanoflagellata), *Monosiga brevicollis* (Choanoflagellata), *Amphimedon queenslandica* (Porifera), *Mnemiopsis leidyi* (Ctenophora), *Nematostella vectensis* (Cnidaria), *Trichoplax adhaerens* (Placozoa), *Drosophila melanogaster* (Arthropoda), *Capitella teleta* (Annelida), *Schmidtea mediterranea* (Platyhelminth), *Caenorhabditis elegans* (Nematoda), as well as *Mus musculus* and *Danio rerio* (Vertebrata) gene models to identify organisms with CPEB orthologs. The cut-off point is set to  $Evalue \leq 0.1$ , because lower Evalues ( $10e-6$ ) would not allow for inclusion of non-orthologous homologs of CPEBs. An additional step to remove redundant sequences is performed by keeping only the longest isoform.

- 2) **Align domains in CPEB homologs:** We will search the putative CPEBs and outgroups identified above with the following HMMs from PFAM: RRM\_1 (PF00076), RRM\_7 (PF16367), and CEBP\_ZZ (PF16366). For CEBP\_ZZ, which should only occur once per protein, we will create an alignment to the CEBP\_ZZ HMM using the `hmm2aln.pl` script which runs `hmmsearch`, `stockholm2fasta`, and custom code to remove indels. For the RRM\_1 and RRM\_7 HMMs, which will often match the same domains, we will run `hmmsearch` separately with each of these HMMs and keep all non-overlapping results, and merge overlapping results by taking the lowest N-terminal coordinate and the highest C-terminal coordinate from the 2 results. We will extract the amino acid sequences between the start and end of the match and use `mafft` with default parameters to generate an alignment of RNA recognition motifs (RRMs).
- 3) Infer phylogenetic relationships among CPEBs and non-CPEB sequences (with similarity to CPEBs) by estimating gene trees of RBD sequence databases from step (3) and step (6).
  - a. IQTREE
  - b. RAxML with 25 starting parsimony trees and 25 random starting trees;
  - c. MrBayes
- 4) Choose best tree for main figure

### 3 WORK COMPLETED SO FAR W DATES

<sampling plan>

<statistical tests>

<inference criteria (p-values, bayes factors, model fit indices)>

<criteria for accepting or rejecting hypotheses>

<description of or link to repo with custom scripts>

<command lines>

### 4 LITERATURE REFERENCES

Paps, J & Holland, PW, 2018, 'Reconstruction of the ancestral metazoan genome reveals an increase in genomic novelty'. *Nature Communications*, vol 9.

Afroz T, Skrisovska L, Belloc E, Guillén-Boixet J, Méndez R, Allain FH. A fly trap mechanism provides sequence-specific RNA recognition by CPEB proteins. *Genes Dev.* 2014;28(13):1498-1514. doi:10.1101/gad.241133.114

## 5 PHYLOTOCOL AMENDMENT HISTORY

| Version | Date          | Significant Revisions  |
|---------|---------------|--|
| 1.1     | Dec. 19, 2020 | Inclusion of E-value cut-off $\leq 0.1$ for step 1   |
| 1.1     | Dec. 19, 2020 | Removal of steps 2, 3, 5 and 6 from previous version (1.0) because non-orthologous sequences will be identified in step 1 under the updated cut-off E-value. |