1  **Exhaustive identification of conserved upstream open reading frames with potential translational**

2  **regulatory functions from animal genomes**

3

4  Hiro Takahashi[1,2#]*, Shido Miyaki[2#], Hitoshi Onouchi[3#], Taichiro Motomura[1], Nobuo Idesako[2], Anna Takahashi[4],

5  Shuichi Fukuyoshi[5], Toshinori Endo[6], Kenji Satou[7] , Satoshi Naito[3,8], and Motoyuki Itoh[9]*

6

7  [1]Graduate School of Medical Sciences, Kanazawa University, Kanazawa 920-1192, Japan

8  [2]Graduate School of Horticulture, Chiba University, Matsudo 271-8510, Japan

9  [3]Graduate School of Agriculture, Hokkaido University, Sapporo 060-8589, Japan

10  [4]Faculty of Information Technologies and Control, Belarusian State University of Informatics and Radio

11  Electronics, Minsk 220013, Belarus

12  [5]Institute of Medical, Pharmaceutical and Health Sciences, Kanazawa University, Kakuma-machi, Kanazawa,

13  Ishikawa 920-1192, Japan

14  [6]Graduate School of Information Science and Technology, Hokkaido University, Sapporo 060-0814, Japan

15  [7]Faculty of Biological Science and Technology, Institute of Science and Engineering, Kanazawa University,

16  Kanazawa 920-1192, Japan

17  [8]Graduate School of Life Science, Hokkaido University, Sapporo 060-0810, Japan

18  [9]Graduate School of Pharmaceutical Science, Chiba University, Chuo-ku, Chiba 260-8675, Japan

19

20  *Correspondence. Tel: +81-76-234-4484; Fax: +81-76-234-4484; Email: takahasi@p.kanazawa-u.ac.jp

21  Correspondence may also be addressed to Motoyuki Itoh. Email: mito@chiba-u.jp

22  [#]Joint first authors.

23

24  **Key words:** upstream open reading frame; translational regulation; bioinformatics; nascent peptide

25

26

27

28    **Abstract**

29    **Background:** Upstream open reading frames (uORFs) are located in the 5′-untranslated regions of many

30    eukaryotic mRNAs, and some peptides encoded in these regions play important regulatory roles in controlling

31    main ORF (mORF) translation. To comprehensively identify uORFs encoding functional peptides, genome-wide

32    searches for uORFs with conserved peptide sequences (CPuORFs) have been conducted in various organisms

33    using comparative genomic approaches. However, in animals, CPuORFs have been identified only by comparing

34    uORF sequences between a limited number of closely related species, and it is unclear how many previously

35    identified CPuORFs encode regulatory peptides.

36    **Results:**  Here, we conducted exhaustive genome-wide searches for animal CPuORFs conserved in various

37    taxonomic ranges, using the ESUCA pipeline, which we recently developed for efficient comprehensive

38    identification of CPuORFs. ESUCA can efficiently compare uORF sequences between an unlimited number of

39    species using BLAST and automatically determine the taxonomic ranges of sequence conservation for each

40    CPuORF. By applying ESUCA to human, chicken, zebrafish, and fruit fly genomes, 1,430 (1,339 novel and 91

41    known) CPuORFs were identified. We examined the effects of 14 human CPuORFs on mORF translation using

42    a transient expression assay. Through this analysis, we identified six novel regulatory CPuORFs that repressed

43    mORF translation in a sequence-dependent manner, all of which were conserved beyond Amniota.

44    **Conclusions:** We discovered a much higher number of animal CPuORFs than previously identified.

45    Furthermore, our results suggest that human CPuORFs conserved beyond Amniota are more likely to encode

46    regulatory peptides than those conserved in narrower taxonomic ranges.

47

## Background

49  The human genome contains many regions encoding potential functional small peptides outside of the

50  well-annotated protein-coding regions (Ingolia, et al., 2014). Some upstream open reading frames (uORFs), which are

51  located in the 5′-untranslated regions (5′-UTRs) of mRNAs, have been shown to encode such functional small

52  peptides. Most uORF-encoded peptides play regulatory roles in controlling the translation of protein-coding main

53  ORFs (mORFs) (Cruz-Vera, et al., 2011; Ito and Chiba, 2013; Morris and Geballe, 2000; Somers, et al., 2013).

54  During the translation of these regulatory uORFs, nascent peptides interact inside the ribosomal exit tunnel to

55  cause ribosome stalling (Bhushan, et al., 2010). Ribosome stalling on a uORF results in translational repression of

56  the downstream mORF because stalled ribosomes block scanning of subsequent pre-initiation complexes and

57  prevent them from reaching the start codon of the mORF (Wang and Sachs, 1997). In some genes, uORF

58  peptides are involved in translational regulation in response to metabolites (Ito and Chiba, 2013).

59  To comprehensively identify uORFs encoding functional peptides, genome-wide searches for uORFs

60  with conserved peptide sequences (CPuORFs) have been conducted using comparative genomic approaches in

61  plants (Hayden and Jorgensen, 2007; Takahashi, et al., 2019; Takahashi, et al., 2012; Tran, et al., 2008; van der

62  Horst, et al., 2018; Vaughn, et al., 2012). To date, 157 CPuORF families have been identified by comparing

63  5′-UTR sequences between plant species. Of these, 101 families were identified in our previous studies by applying

64  our original methods, BAIUCAS (Takahashi, et al., 2012) and ESUCA (an advanced version of BAIUCAS)

65  (Takahashi, et al., 2019) to genomes of *Arabidopsis*, rice, tomato, poplar, and grape.

66  ESUCA has many unique functions (Takahashi, et al., 2019), such as efficient comparison of uORF

67  sequences between an unlimited number of species using BLAST, automatic determination of taxonomic ranges

68  of CPuORF sequence conservation, systematic calculation of $K_a$/$K_s$ ratios of CPuORF sequences, and wide

69  compatibility with any eukaryotic genome whose sequence database is registered in ENSEMBL (Zerbino, et al.,

70  2018). More importantly, to distinguish between 'spurious' CPuORFs conserved because they encode parts of

71  mORF-encoded proteins and 'true' CPuORFs conserved because of the functions of their encoded small peptides,

72  ESUCA assesses whether a transcript containing a fusion of a uORF and an mORF is a major or minor form

73  among homologous transcripts (Takahashi, et al., 2019). By using these functions, ESUCA can efficiently

74  identify CPuORFs likely to encode functional small peptides. In fact, our recent study demonstrated that poplar

75  CPuORFs encoding regulatory peptides were efficiently identified by selecting ones conserved across diverse

76  eudicots using ESUCA (Takahashi, et al., 2019).

77        To date, only a few studies on genome-wide identification of animal CPuORFs have reported. In these

78  previous studies, uORF sequences were compared between a limited number of closely related species, such as

79  human and mouse or several species in dipteran, leading to identification of 204 and 198 CPuORFs in human

80  and mouse, respectively (Crowe, et al., 2006), and 44 CPuORFs in fruit fly (Hayden and Bosco, 2008).

81  Additionally, the relationships between taxonomic ranges of CPuORF conservation and the likelihood of having

82  a regulatory function have not been studied in animals.

83        Accordingly, in this study, we applied ESUCA to genomes of fruit fly, zebrafish, chiken, and human to

84  exhaustively identify animal CPuORFs and to determine the taxonomic range of their sequence conservation.

85  Using ESUCA, we identified 1,430 animal (1,339 novel and 91 known) CPuORFs belonging to 1,337 CPuORF

86  families. We examined the effects of 15 CPuORFs conserved in various taxonomic ranges on mORF translation,

87  using a transient expression assay. Through this analysis, we identified six novel regulatory CPuORFs that

88  repress mORF translation in a sequence-dependent manner. All of the six regulatory CPuORFs are conserved

89  beyond Amniota, suggesting that human CPuORFs conserved beyond Amniota are more likely to encode

90  functional peptides than those conserved in narrower taxonomic ranges.

91

92

93  ## Materials and methods

94  **Extraction of CPuORFs using ESUCA**

95  ESUCA was developed as an advanced version of BAIUCAS (Takahashi et al., 2012) in our previous study

96  (Takahashi et al., 2019). ESUCA consists of six steps, and some of these steps are divided into substeps, as

97  shown in Fig. 1A and 1B. To identify animal CPuORFs using ESUCA, the following eight-step procedures were

98  conducted, including the six ESUCA steps: 0) data preparation for ESUCA, 1) uORF extraction from the 5′-UTR,

99   2) calculation of uORF-mORF fusion ratios, 3) uORF-tBLASTn against transcript sequence databases, 4)

100  mORF-tBLASTn against downstream sequence datasets for each uORF, 5) calculation of $K_a$/$K_s$ ratios, 6)

101  determination of the taxonomic range of uORF sequence conservation, and 7) manual validation after ESUCA.

102  See the Supplementary Materials and Methods for details.

103

104  **Determination of the taxonomic range of uORF sequence conservation for animal CPuORFs**

105  To apply ESUCA to animal genomes, we defined 19 animal taxonomic categories, as shown in Fig. 1C. See the

106  Supplementary Materials and Methods for details.

107

108  **Plasmids and reporter assays**

109  DNA fragments containing a control CPuORF (Con) or the frameshift mutant version (fs) of the 15 selected

110  genes were subcloned into pSV40:UTR:Fluc. Reporter assays were conducted using a Dual-Luciferase Reporter

111  Assay system (Promega, Madison, WI, USA). See the Supplementary Materials and Methods for details.

112

113  # Results

114  **Genome-wide search for animal CPuORFs using ESUCA**

115  Prior to ESUCA application (Fig. 1A and 1B), we counted the number of protein-coding genes for four species,

116  i.e., fruit fly, zebrafish, chiken, and human. As shown in Supplementary Table S1, 13,938, 25,206, 14,697, and

117  19,956 genes were extracted for fruit fly, zebrafish, chiken, and human, respectively. After step 1 of ESUCA, we

118  calculated the numbers of uORFs and protein-coding genes with any uORF for each species. As shown in

119  Supplementary Table S1, 17,035 (7,066), 39,616 (14,453), 8,929 (3,535), and 44,085 (12,321) uORFs (genes)

120  were extracted from fruit fly, zebrafish, chicken, and human genomes, respectively. In this analysis, when

121  multiple uORFs from a gene shared the same stop or start codon, they were counted as one. Potential candidate

122  CPuORFs were narrowed down by selection at step 2.5 of ESUCA in a step-by-step manner, as shown in

123  Supplementary Table S1. The numbers of BLAST hits (expressed sequence tag [EST], transcriptome shotgun

124  assembly [TSA], assembled EST/TSA, and RefSeq RNA sequences) extracted at step 3.5 are also shown in

125     Supplementary Table S1. After the final step of ESUCA, 49, 195, 235, and 1,453 candidate CPuORFs were

126     extracted from fruit fly, zebrafish, chiken, and human, respectively. We conducted manual validation for the

127     extracted candidate CPuORFs as described in our previous study (Takahashi, et al., 2019) and in the

128     Supplementary Materials and Methods. We selected CPuORFs conserved in at least two orders other than the

129     order to which the original species belongs. In total, 1,430 animal CPuORFs (35 for fruit fly, 151 for zebrafish,

130     206 for chicken, and 1,038 for human) were identified (Fig. 1D). Of these, 1,339 CPuORFs were newly

131     identified in the current study. Detailed information on the identified CPuORFs is shown in Supplementary Table

132     S2. The identified CPuORF-containing genes were classified into 1,124 ortholog groups on the basis of

133     similarities of mORF-encoded amino acid sequences, using OrthoFinder ver. 1.1.4 (Emms and Kelly, 2015).

134     CPuORFs with similar amino acid sequences from the same ortholog groups were categorized as the same

135     CPuORF families (homology groups [HGs]; Supplementary Materials and Methods). The identified 1,430

136     CPuORFs were classified into 1,337 HGs. We assigned HG numbers to 1,337 HGs in an order based on numbers

137     of orders in which any CPuORF belonging to each HG was extracted, the taxonomic range of the sequence

138     conservation of each HG, and gene ID numbers. When multiple CPuORF families were identified in the same

139     ortholog groups, the same HG number with a different subnumber was assigned to each of the families.

140

141     **Sequence-dependent effects of CPuORFs on mORF translation**

142     To address the relationship between taxonomic ranges of CPuORF conservation and likelihood of having

143     regulatory function, we selected 15 human CPuORFs conserved in various taxonomic ranges, including a

144     previously identified sequence-dependent regulatory CPuORF, the *PTP4A1* CPuORF (Hardy, et al., 2019), as a

145     positive control, and examined their sequence-dependent effects on the expression of the downstream reporter

146     gene using transient expression assays (Fig. 2). Other uORFs overlapping any of the selected CPuORFs were

147     eliminated by introducing mutations that changed the ATG codons of the overlapping uORFs to other codons but

148     did not alter the amino acid sequences of the CPuORFs (Supplementary Fig. S5). The resulting modified

149     CPuORFs were used as Con CPuORFs (Fig. 2B). To assess the importance of amino acid sequences for the

150     effects of these CPuORFs on mORF translation, fs mutations were introduced into the Con CPuORFs such that

151  the amino acid sequences of their conserved regions could be altered (see Supplementary Materials and Methods

152  and Supplementary Fig. S1 for details). In seven of the 15 CPuORFs, the introduced frameshift mutations

153  significantly upregulated the expression of the reporter gene, indicating that these CPuORFs repressed mORF

154  translation in a sequence-dependent manner (Fig. 2C). All of the seven CPuORFs with the sequence-dependent

155  repressive effects were conserved beyond Amniota (Fig. 2A). In contrast, any of four CPuORFs conserved only

156  among Amniota did not show significant sequence-dependent effects (Fig. 2C). These results suggest that human

157  CPuORFs conserved beyond Amniota are more likely to encode regulatory peptides than those conserved in

158  narrower taxonomic ranges.

159

## Discussion

161  In the current study, by applying ESUCA to four animal genomes, we identified 1,430 CPuORFs belonging to

162  1,337 HGs. Taxonomic ranges of sequence conservation of these CPuORFs largely vary, demonstrating that

163  ESUCA can identify CPuORFs conserved in various taxonomic ranges (Supplementary Table S3). Moreover,

164  seven of 11 human CPuORFs conserved beyond Amniota exhibited sequence-dependent repressive effects on

165  mORF translation, whereas all four CPuORFs conserved only among Amniota showed no significant

166  sequence-dependent effects. This result suggest that human CPuORFs conserved beyond Amniota are more

167  likely to encode regulatory peptides than those conserved in narrower taxonomic ranges. Of the 1,038 CPuORFs

168  extracted from the human genome, 78 are conserved beyond Amniota (Supplementary Table S3). Therefore,

169  these 78 CPuORFs are promising candidates of regulatory CPuORFs encoding peptides that control mORF

170  translation. CPuORFs encoding functional peptides may also be found among the remaining human CPuORFs

171  conserved in narrower taxonomic ranges because the $K_a/K_s$ ratios suggest that purifying selection acted on the

172  amino acid sequences of these CPuORFs.

173        In this study, we identified six novel human regulatory CPuORFs (in the *MKKS*, *SLC6A8*, *FAM13B*,

174  *MIEF1*, *KAT6A*, and *LRRC8B* genes) with sequence-dependent repressive effects on mORF translation. Of these,

175  the *MKKS* CPuORF is a translational regulator that represses the production of a protein involved in

176  McKusick-Kaufman syndrome (Akimoto, et al., 2013); however, the amino acid sequence dependence of the

177 CPuORF function was not reported. Interestingly, the *MIEF1* CPuORF-encoded peptide is a functional peptide

178 localized in the mitochondria (Samandi, et al., 2017). Thus, the *MIEF1* CPuORF may have dual functions.

179 Chemical screening recently identified a compound that causes nascent peptide-mediated ribosome

180 stalling in the mORF of the human *PCSK9* gene, resulting in specific translational inhibition of *PCSK9* and a

181 reduction in total plasma cholesterol levels (Lintner, et al., 2017). Nascent peptide-mediated ribosome stalling in

182 some of the previously identified regulatory CPuORFs is promoted by metabolites, such as polyamine, arginine,

183 and sucrose (Ito and Chiba, 2013; Yamashita, et al., 2017). Therefore, compounds that promote nascent

184 peptide-mediated ribosome stalling in CPuORFs could be identified by chemical screening through a method

185 similar to that used for the screening of the stall-inducing compound for *PCSK9*. The data from the current study

186 may be useful for selection of CPuORFs as potential targets for pharmaceutical drugs and for identification of

187 regulatory CPuORFs.

188

193

194 **Competing interests**: none declared.

195

196

197

198

199

200

201

202

203

# References

204

205    Akimoto, C*., et al.* (2013) Translational repression of the McKusick-Kaufman syndrome transcript by unique

206        upstream open reading frames encoding mitochondrial proteins with alternative polyadenylation sites,

207        *Biochim Biophys Acta*, **1830**, 2728-2738.

208    Bhushan, S*., et al.* (2010) Structural basis for translational stalling by human cytomegalovirus and fungal arginine

209        attenuator peptide, *Mol. Cell*, **40**, 138-146.

210    Crowe, M.L., Wang, X.Q. and Rothnagel, J.A. (2006) Evidence for conservation and selection of upstream open

211        reading frames suggests probable encoding of bioactive peptides, *BMC Genomics*, **7**, 16.

212    Cruz-Vera, L.R*., et al.* (2011) Nascent polypeptide sequences that influence ribosome function, *Curr. Opin.*

213        *Microbiol.*, **14**, 160-166.

214    Emms, D.M. and Kelly, S. (2015) OrthoFinder: solving fundamental biases in whole genome comparisons

215        dramatically improves orthogroup inference accuracy, *Genome Biol.*, **16**, 157.

216    Hardy, S*., et al.* (2019) Magnesium-sensitive upstream ORF controls PRL phosphatase expression to mediate

217        energy metabolism, *Proc Natl Acad Sci U S A*, **116**, 2925-2934.

218    Hayden, C.A. and Bosco, G. (2008) Comparative genomic analysis of novel conserved peptide upstream open

219        reading frames in *Drosophila melanogaster* and other dipteran species, *BMC Genomics*, **9**, 61.

220    Hayden, C.A. and Jorgensen, R.A. (2007) Identification of novel conserved peptide uORF homology groups in

221        *Arabidopsis* and rice reveals ancient eukaryotic origin of select groups and preferential association with

222        transcription factor-encoding genes, *BMC Biol.*, **5**, 32.

223    Ingolia, N.T*., et al.* (2014) Ribosome profiling reveals pervasive translation outside of annotated protein-coding

224        genes, *Cell Rep.*, **8**, 1365-1379.

225    Ito, K. and Chiba, S. (2013) Arrest peptides: cis-acting modulators of translation, *Annu. Rev. Biochem.*, **82**,

226        171-202.

227    Lintner, N.G*., et al.* (2017) Selective stalling of human translation through small-molecule engagement of the

228      ribosome nascent chain, *PLoS Biol*, **15**, e2001882.

229      Morris, D.R. and Geballe, A.P. (2000) Upstream open reading frames as regulators of mRNA translation, *Mol.*

230           *Cell Biol.*, **20**, 8635-8642.

231      Samandi, S., *et al.* (2017) Deep transcriptome annotation enables the discovery and functional characterization of

232           cryptic small proteins, *Elife*, **6**.

233      Somers, J., Poyry, T. and Willis, A.E. (2013) A perspective on mammalian upstream open reading frame

234           function, *Int. J. Biochem. Cell Biol.*, **45**, 1690-1700.

235      Takahashi, H., *et al.* (2019) ESUCA: a pipeline for genome-wide identification of upstream open reading frames

236           with evolutionarily conserved sequences and determination of the taxonomic range of their conservation,

237           *bioRxiv*, 524090.

238      Takahashi, H., *et al.* (2012) BAIUCAS: a novel BLAST-based algorithm for the identification of upstream open

239           reading frames with conserved amino acid sequences and its application to the *Arabidopsis thaliana* genome,

240           *Bioinformatics*, **28**, 2231-2241.

241      Tran, M.K., Schultz, C.J. and Baumann, U. (2008) Conserved upstream open reading frames in higher plants,

242           *BMC Genomics*, **9**, 361.

243      van der Horst, S., *et al.* (2018) Novel pipeline identifies new upstream ORFs and non-AUG initiating main ORFs

244           with conserved amino acid sequences in the 5' leader of mRNAs in *Arabidopsis thaliana*, *RNA*, **25**,

245           292-304.

246      Vaughn, J.N., *et al.* (2012) Known and novel post-transcriptional regulatory sequences are conserved across plant

247           families, *RNA*, **18**, 368-384.

248      Wang, Z. and Sachs, M.S. (1997) Ribosome stalling is responsible for arginine-specific translational attenuation

249           in *Neurospora crassa*, *Mol. Cell Biol.*, **17**, 4904-4913.

250      Yamashita, Y., *et al.* (2017) Sucrose sensing through nascent peptide-meditated ribosome stalling at the stop

251           codon of Arabidopsis *bZIP11* uORF2, *FEBS Lett.*, **591**, 1266-1277.

252      Zerbino, D.R., *et al.* (2018) Ensembl 2018, *Nucleic Acids Res.*, **46**, D754-D761.

253

254

## **Figure Legends**

256 **Figure 1.** Identification of animal CPuORFs using ESUCA. (A) Data preparation. (B) Outline of the ESUCA

257 pipeline. Numbers with parenthesis indicate datasets labeled with the same numbers in A. (C) Defined animal

258 taxonomic categories. (D) Numbers of identified CPuORFs.

259

260 **Figure 2.** Taxonomic conservation and experimental validation of 15 selected human CPuORFs. (A) Taxonomic

261 ranges of conservation of CPuORFs examined in transient assays. Filled cells in each taxonomic category

262 indicate the presence of uORF-tBLASTn and mORF-tBLASTn hits for CPuORFs of the indicated genes. (B)

263 Reporter constructs used for transient assays. The hatched box in the frameshift (fs) mutant CPuORF indicates

264 the frame-shifted region. Dotted boxes represent the first five nucleotides of the mORFs associated with the 15

265 human CPuORFs. (C) Relative luciferase activities of control (white) or frameshift (gray) CPuORF reporter

266 plasmids. Means ± SDs of three biological replicates are shown. $*p < 0.05$.

267

268

269

270

271

272

273

274

275

276

277

278

279

- 11 -

280

## Supplementary Figure Legend

282 **Figure S1.** Nucleotide sequences of the 5′-UTRs and amino acid sequences of the CPuORFs analyzed in this

283 study. (A–O) The 5′-UTRs of *PTP4A1* (A), *MKKS* (B), *SLC6A8* (C), *FAM13B* (D), *MIEF1* (E), *EIF5* (F),

284 *MAPK6* (G), *MEIS2* (H), *KAT6A* (I), *SLC35A4* (J), *LRRC8B* (K), *CDH11* (L), *PNRC2* (M), *BACH2* (N), and

285 *FGF9* (O). The nucleotide sequences of the CPuORFs are shown in bold. The deduced amino acid sequences of

286 the control (Con) and frameshift (fs) CPuORFs are indicated. The nucleotide sequences of other uORFs are

287 underlined with a bold line. Dotted underlines indicate the nucleotide sequences of other uORFs overlapping the

288 CPuORFs, whose initiation codons were altered to other codons by introducing nucleotide substitutions that did

289 not change the amino acid sequences of the CPuORFs. The replaced nucleotides are shown as white letters in a

290 black background. The nucleotides that were deleted and inserted in the frameshift mutants are shaded. The main

291 coding sequences that were contained in the reporter constructs are boxed. The shaded boxes indicate the

292 nucleotides changed to avoid the appearance of in-frame termination codons. Cloning sites were added at either

293 end of the nucleotide sequences in controls to be subcloned into plasmid pGL4.10 with an SV40 promoter

294 (pSV40:5′UTR::luc2) are underlined (Fig. 2B).

295

# Figure 1



**A** — Data preparation before ESUCA

**Step 0.1** Transcript dataset construction based on genome information for a certain organism

Genome sequences (FASTA) data and genomic coordinates (GFF3)

(1a) 5'UTR — mORF RNA → Translation → (1b) mORF Protein

**Step 0.2** Transcript base sequence dataset construction from EST/TSA/RefSeq RNA

EST+TSA+RefSeq sequences with taxonomy DB — Extraction → (5) Taxonomy DB

Extraction → (2) RefSeq — Assembling & mapping → (3) Assembled EST/TSA — Extraction → (4) EST/TSA/RefSeq

**B** — ESUCA pipeline

**Step 1** Extraction of uORF sequences from 5'UTR ← (1a)

**Step 2** Calculation of uORF-mORF fusion ratios ← (2)

uORF-mORF fusion ratio ≥ 0.3 → Discarded

**Step 3.1** Homology searches of the uORF amino acid sequences against transcript sequence database (uORF-tBLASTn) ← (2)+(3)

**Step 3.2** Extraction of uORF from each uORF-tBLASTn hit sequence

**Step 4.0** Extraction of downstream sequences of the ORFs matching the original uORF → Downstream sequence dataset for each uORF

**Step 4.1** Selection of uORFs conserved between homologous genes (mORF-tBLASTn) ← (1b)

**Step 4.2** Removal of contaminant ESTs and TSAs (BLASTn) ← (4)

**Step 4.3** Selection of uORFs conserved between homologs from more than two orders ← (5)

**Step 5** Calculation of $K_a/K_s$ ratios ← Representative uORF from each order

$K_a/K_s$ ratio ≥ 0.5  or $q$ value ≥ 0.05 → Discarded

**Step 6** Determination of the taxonomic range of uORF sequence conservation ← (5)

**Step 7** Manual validation after ESUCA

**C**

Metazoa
Eumetazoa
Cnidaria
Bilateria
Protostomia
Ecdysozoa
Arthropoda
Insecta

Deuterostomia
Chordata
Vertebrata
Euteleostomi
Actinopterygii
Ostarioclupeomorpha

Sarcopterygii
Tetrapoda
Amphibia
Amniota
Sauropsida
Aves

Mammalia
Eutheria
Euarchontoglires

**D**

| Species | CPuORF |
|---|---|
| *D. melanogaster* | 35 |
| *D. rerio* | 151 |
| *G. gallus* | 206 |
| *H. sapiens* | 1038 |

Number of animal CPuORFs: 1430
Number of animal HGs: 1337

# Figure 2

**A**

| | | | | | PTP4A1 | MKKS | SLC6A8 | FAM13B | MIEF1 | EIF5 | MAPK6 | MEIS2 | KAT6A | SLC35A4 | LRRC8B | CDH11 | PNRC2 | BACH2 | FGF9 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Metazoa | | | | | Cnidaria | ■ | | | | | | | | | | | | | | |
| | Bilateria | | | | Protostomia | ■ | ■ | | | | | | | | | | | | | |
| | | Deuterostomia | | | Deuterostomia | ■ | ■ | | | | | | | | | | | | | |
| | | | Vertebrata | | Vertebrata | ■ | | ■ | ■ | ■ | ■ | ■ | ■ | ■ | | | | | | |
| | | | | | Actinopterygii | ■ | | ■ | ■ | ■ | ■ | ■ | ■ | | ■ | | | | | |
| | | | Tetrapoda | | Amphibia | ■ | | ■ | ■ | ■ | ■ | ■ | | | ■ | ■ | | | | |
| | | | | Amniota | Sauropsida | ■ | | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ |
| | | | | | Mammalia | ■ | | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ | ■ |

**B**

*SV40::UTR(Con):Fluc*

SV40pro — CPuORF (Con) — *Fluc* — ter

*SV40::UTR(fs):Fluc*

SV40pro — CPuORF (fs) — *Fluc* — ter

**C**



Relative Fluc activity (Con=1)