

Reviewers' Comments:

Reviewer #1:

Remarks to the Author:

In this manuscript, Zhang et al describe a large scale evolutionary analysis of uORFs across multiple eukaryotic species. Even though the conclusions made based on this analysis are in a general agreement with what we already know about uORFs, the study is unprecedented in its scale and is, therefore, of general interest. Given the intense attention that the topic of small ORFs translation received recently, the manuscript is also timely.

Nonetheless, I found the manuscript to be insufficiently clear in several places. Specifically, I have the following comments.

General:

1. The manuscript gives the impression that it investigates uORFs across all eukaryotes. This, however, is not true. While the study is unprecedented in its scale it is limited to multicellular eukaryotes such as plants and animals. Large phylogenetic clusters, such as fungi and protists are not represented in this study. I strongly suspect that the evolution, distribution and function of uORFs may significantly differ in the organisms from these phyla. Just to give an example, consider recently discovered genetic codes in some ciliates where termination of translation takes place only in close proximity to mRNA 3' ends, e.g. in *Condylostoma magnum* all stop codons code for amino acids in internal positions of mRNA and in *Euplotes* stop codons cause +1 or +2 frameshifting unless in close proximity to the 3' end. In these organisms, once ribosomes initiate translation they are expected to continue the translation of the entire mRNA. Thus, short uORFs are impossible. Given the phylogenetic diversity of many protists and our very limited knowledge of their molecular biology, we may expect many surprising findings regarding the organization of their genetic information and considerable differences from what has been revealed in this manuscript. Thus, at a minimum, the authors should clearly define the phylogenetic boundaries of their study, e.g. "uORFs in plants and animals" instead of "uORFs in eukaryotes", but perhaps it would also be good if the authors discuss the potential limitations of extrapolating their findings on the entire eukaryotic kingdom.

2. One of the authors' conclusions is that most uORFs are regulatory rather than coding for functional peptides, while correctly acknowledging that some uORFs do code for functional peptides. While such a statement is most likely true, it is also vague and hence not very informative. First, "most" stands for "more than a half" which could be 51% or 99%. I wonder if the authors could try to give a more quantitative estimate. In doing so, I also suggest that the authors should take care in defining what they consider functional or perhaps even avoiding the use of the term 'functional', so not to get into a type of controversy such as the one that took place when ENCODE claimed that 80% of the human genome is functional. It seems to me that by function here the authors mean evidence of evolutionary selection. Not all functions are under evolutionary selection, consider human olfactory receptors, many of which, although clearly functional, do accumulate deleterious mutations and evolve almost neutrally. At the same time not all uORFs that exhibit $\omega \ll 1$ necessarily encode functional peptides, because some uORFs are known to alter ribosome movement by making ribosomes stall via specific interactions inside the peptide channel. Such stalling peptides may not function on their own outside of the ribosome even though they would be expected to evolve as protein-coding.

Specific

1. It is not clear how exactly the groups of genes were divided into the categories for the analyses shown in Fig. 1Sb. A more detailed explicit description is necessary.

2. The authors extensively used the data from ref. 27 (McGillivray et al) and attempted to make certain conclusions regarding the evolution of uORFs reported in that work, for example, they found the evidence that these uORFs are more conserved. This is inappropriate. McGillivray et al used conservation as one of the features used for their machine learning algorithm: "Features were chosen to cover a broad range of categories of data, including features associated with uORF position and length, conservation, functional metrics like RNA expression, and sequence-based signatures that may relate to translation." It makes no sense to show increased conservation of uORFs that were predicted based on their conservation.

3. Comparing the strength of Kozak context in uORFs and CDS ATGs. If we take two groups of sequences and compare, they are likely to differ in some respect. If we define one as optimal, the other would become suboptimal. Therefore, the purpose of this analysis is unclear to me. Perhaps it would be more meaningful to compare three groups of contexts rather than two, by adding ATGs that are not used for initiation, e.g. internal ATGs from CDS or ATGs from 3' UTRs or intergenic regions. We expect that the context of such ATGs should not evolve to optimize translation initiation and would provide an estimate for a background context and a variation in contexts. Then we would expect that uORFs context should be optimized for translation initiation, but not as strong as CDS ATG. By having three points the authors could estimate whether uORF ATG context is closer to neutral or that of CDS.

4. "Unsurprisingly, for both uORFs and CDSs, the distance between two species from a clade tended to be significantly shorter than that between one species in that clade and another species outside of that clade (Fig. 6c). These results suggest that the Kozak contextual characteristics tend to be similar between closely related species for both uORFs and CDSs."

This is indeed so unsurprising that it is unclear why was it even done. I believe that any other sequence, e.g. a context of stop codons would exhibit the same behaviour.

5. The authors made an observation that uORFs occurrence anticorrelates with expression levels. This makes sense, but there could be at least two reasons for that. One is that the regulation usually works by suppression, hence the mRNAs whose translation is regulated by uORFs are likely to be lowly expressed. The other is that the negative selection acting on uAUGs is expected to be weaker for lowly expressed mRNAs. These two scenarios are drastically different, could authors try to estimate contributions of each of these two scenarios?

6. To demonstrate the evidence of positive selection on 162 newly fixed uORFs, authors have used the asymptotic McDonald-Kreitman, where the alpha parameter is the proportion of substitutions that are due to adaptive evolution. But confidence intervals are quite wide and contain zeros (as well as negative values), so there seems to be no strong evidence of positive selection (Fig 2a)

7. It is unclear how the relative fixation probability of newly originated uORFs was calculated. Could the authors provide an explicit description of the procedure?

Reviewer #2:

Remarks to the Author:

In the study of Zhang et al., the authors analyzed more than 10 million "uORFs" in over 200 eukaryotic species. They found that 1) most of "uORFs" are under purifying selection. 2) the coding region of "uORFs" is overall less conserved, suggesting that uORF is under neutral evolution or weak selective pressure. Finally, they also analyzed the evolution of start codon and flanking context of uORFs. While the manuscript is written well, many of main conclusions are not new, which have been reported by previous studies. Although previous studies analyzed uORF evolution usually based on a small subset of closely-related species, simply using more species does not

significant extend our knowledge on the origin of uORF translation and its evolution. My major concerns are as follows.

1) Like canonical translation, uORF translation is energy-consuming. Uncontrolled uORF translation may inhibit translation in main CDS. Therefore, it is not unexpected that the potential of uORF translation in 5' UTR has been eliminated during evolution. Also, similar observations have been reported by previous studies, for example PMC5793785, PMC4890304.

2) uORF translation plays various roles in gene expression regulations. As demonstrated by many previous studies (see reviews PMID: 28698598, PMID: 31003826), uORF may encode functional peptide, or uORF translation may control downstream translation in main CDS. Again, it is not unexpected that the coding region of uORFs may not under negative selection, if they do not encode functional peptides.

3) About uORF definition. In this study, uORF is defined as a 5' UTR region starting with ATG and ending with an in-frame stop codon (TAG, TAA or TGA). The uORF definition is problematic. First, they overlooked uORFs starting with non-canonical start codons such as CTG, TTG or ATT. Previous studies have suggested that non-canonical start codons are more prevalence than canonical start codon (i.e. ATG) in uORFs. Second, I would define these regions as putative uORFs or potential uORFs, because majority of these so-called "uORFs" are not translatable. Only a very small number of these putative uORFs are real uORFs with significant protein translation. Analysis based on these putative uORFs will be strongly affected by huge amount of false positives (or background noises), and can not be used to support the conclusions on uORFs.

For example,

i) they found O/E ratio (based on these putative uORFs) is significant lower than 1, suggesting that "purifying selection is the major force shaping the prevalence of uORFs". This result only suggests that ATG triplets are depleted in 5' UTR. Purifying selection for ATG triplets in 5' UTR does not mean a necessary of selection for uORFs. In fact, at least in yeast, a previous study (PMC5793785) reported an elevated non-canonical start codon in 5' UTR, indicating a possibility to maintain some kinds of uORF translation.

ii) The authors found that the dN/dS ratio for uORF CDS is "roughly equal to 1 between human and macaque". They concluded that this result supports neutral evolution of uORFs. However, because majority of "uORFs" in their datasets are non-translatable (or not real uORF), these negative uORFs may significantly increase the dN/dS ratio, since they encode nothing. Again, in *Drosophila*, the dN/dS ratio for all uORF CDSs is close to 1, but later, they found that "uORFs with higher Kozak scores presented significantly lower dN/dS ratio in *Drosophila*, suggesting a scenario in which the coding regions of uORFs with optimal Kozak sequence context are under stronger purifying selection in *Drosophila*". Since ATG surrounded by Kozak sequences are more likely to be translated, I believe their negative result (i.e. dN/dS is close to 1) is due to too many negative uORFs in their datasets.

iii) the same problem can be found in the analysis of "evolution of contextual characteristics that influence uORF translation".

4) Page 4, line 94. "gene expression level is a major determinant of the uORF distribution across genes in a eukaryotic species" Because the number of putative uORFs positively correlates with 5' UTR length, I wondered whether 5' UTR may confound the correlation of putative uORF number to gene expression.

5) Page 4, line 101, "while maintaining the same dinucleotide frequency". Please explain why dinucleotide frequency is maintained. Does single, or trip-nucleotide frequency significantly affect O/E ratio?

6) O/E ratio in 5' UTR might be ok to estimate the selection for ATG triplets. To strength the results, O/E ratio in 3' UTR should be considered as negative control, since translation in 3' UTR ORFs is less likely than that in 5' UTR. In addition, it would be great if O/E ratios for the other 61 triplets are displayed.

7) Page 5, line 122. "The O/E ratio varied wildly across the 216 species", is this ratio affected by different background ATG frequency (E) across the species?

8) Page 8, they found longer uORFs have fewer conserved peptides. This is a little unexpected to me. Because uORF translation is energy-consuming. If a uORF plays regulator role, a shorter ORF is sufficient to block ribosome scanning to downstream region. The longer ORF does not significant

benefit the regulator role, but indeed consume more energy.

1 **Response to Reviewers' comments**

2 **Reviewer #1 (Remarks to the Author):**

3 In this manuscript, Zhang et al describe a large scale evolutionary analysis of uORFs across multiple eukaryotic
4 species. Even though the conclusions made based on this analysis are in a general agreement with what we
5 already know about uORFs, the study is unprecedented in its scale and is, therefore, of general interest. Given
6 the intense attention that the topic of small ORFs translation received recently, the manuscript is also timely.
7 Nonetheless, I found the manuscript to be insufficiently clear in several places. Specifically, I have the following
8 comments.

9 **Response:** We greatly appreciate the enthusiasm and the positive feedback from this reviewer. The comments
10 and suggestions are precious and very helpful for us to make this revision. In this revision, we have fully
11 considered your comments and made the revisions accordingly. Please refer to the point-to-point response for
12 details.

13

14 General:

15 1. The manuscript gives the impression that it investigates uORFs across all eukaryotes. This, however, is not
16 true. While the study is unprecedented in its scale it is limited to multicellular eukaryotes such as plants and
17 animals. Large phylogenetic clusters, such as fungi and protists are not represented in this study. I strongly
18 suspect that the evolution, distribution and function of uORFs may significantly differ in the organisms from
19 these phyla. Just to give an example, consider recently discovered genetic codes in some ciliates where
20 termination of translation takes place only in close proximity to mRNA 3' ends, e.g. in *Condylostoma magnum*
21 all stop codons code for amino acids in internal positions of mRNA and in *Euplotes* stop codons cause +1 or +2
22 frameshifting unless in close proximity to the 3' end. In these organisms, once ribosomes initiate translation
23 they are expected to continue the translation of the entire mRNA. Thus, short uORFs are impossible. Given the
24 phylogenetic diversity of many protists and our very limited knowledge of their molecular biology, we may
25 expect many surprising findings regarding the organization of their genetic information and considerable
26 differences from what has been revealed in this manuscript. Thus, at a minimum, the authors should clearly
27 define the phylogenetic boundaries of their study, e.g. "uORFs in plants and animals" instead of "uORFs in
28 eukaryotes", but perhaps it would also be good if the authors discuss the potential limitations of extrapolating
29 their findings on the entire eukaryotic kingdom.

30 **Response:** We thank this reviewer for pointing out this issue. In this revision, we added 242 fungi and 23 protists
31 into our analysis to cover all the large phylogenetic clusters of eukaryotes. The annotated putative canonical
32 uORFs in fungi and protists are summarized as follows (Page 4, Lines 78-84):

33 "The number of annotated protein-coding genes in the 242 fungus genomes ranged from 3,623
34 (*Pneumocystis murina*) to 32,847 (*Fibularhizoctonia sp.*). We identified a total of 3,469,095 uORFs in these
35 fungal genomes, with the number of uORFs ranging from 1,233 (*Malassezia sympodialis*) to 94,695
36 (*Verticillium longisporum*) (Supplementary Table 1). Among the 23 protists, the number of annotated protein-
37 coding genes ranged from 3,398 (*Condylostoma magnum*) to 38,544 (*Emiliana huxleyi*), and the number of

38 uORFs ranging from 1,903 (*Plasmodium falciparum*) to 99,859 (*Cystoisospora suis*), which resulted in a total
39 of 434,267 uORFs in these protist genomes (Supplementary Table 1).”

40 The downstream analyses including the O/E ratio comparisons (Supplementary Figs. 2-4), the influence of gene
41 expression on uORF occurrences (Supplementary Fig. 5), selective constraints on uORF sequences
42 (Supplementary Fig. 8 and 10), and the start codon sequence context (Fig. 6; Supplementary Fig. 13 and 16)
43 was performed for these newly included species as well. Please refer to the updated analyses in TEXT.

44 Although the patterns of uORF distribution and sequence evolution in fungi and protists are largely
45 consistent with that in multicellular animals and plants, some differences were indeed observed in several
46 protists. However, our analysis suggests that the overall occurrence of uORFs are still under purifying selection
47 in these species (Page 14, Lines 453-466):

48 “Among the 481 eukaryotes we studied, the O/E ratio of uORFs was significantly less than 1 in all the 216
49 multi-cellular and the 242 fungal species. Such a pattern was observed in only 17 of the 23 protists, however,
50 and the O/E ratio of uORFs was close to or higher than 1 in the remaining six protists, including *Condylostoma*
51 *magnum* (1.041, 95% CI 1.031~1.050), *Cystoisospora suis* (1.161, 95% CI 1.154~1.169), *Toxoplasma gondii*
52 (0.998, 95% CI 0.989~1.1.007), *Nannochloropsis gaditana* (0.997, 95% CI 0.986~1.007), and two malaria
53 vectors *Plasmodium yoelii* (1.016, 95% CI 1.008~1.025) and *Plasmodium vivax* (0.989, 95% CI 0.975~1.004).
54 It is well established that in protists the nuclear genetic code frequently changed, mostly due to stop codon
55 reassignments¹. In particular, *C. magnum* has no dedicated stop codons², and every uORF is supposed to
56 terminate near the end of a transcript and overlaps with the main CDS. Interestingly, the O/E ratio of uORFs in
57 the 5' UTR regions that are proximal to CDS (within 100 nt or 150 nt) were significantly lower than 1 in five of
58 the six protists except *C. magnum* (Supplementary Table 11). Thus, our results suggest that overall uORF
59 occurrence in 5' UTRs of protists is still under purifying selection, however, whether and how the genetic code
60 reassignments affect the distribution and evolution of uORFs in certain protists deserve further study.”

61

62 2. One of the authors' conclusions is that most uORFs are regulatory rather than coding for functional peptides,
63 while correctly acknowledging that some uORFs do code for functional peptides. While such a statement is
64 most likely true, it is also vague and hence not very informative. First, “most” stands for “more than a half”
65 which could be 51% or 99%. I wonder if the authors could try to give a more quantitative estimate. In doing so,
66 I also suggest that the authors should take care in defining what they consider functional or perhaps even
67 avoiding the use of the term ‘functional’, so not to get into a type of controversy such as the one that took place
68 when ENCODE claimed that 80% of the human genome is functional. It seems to me that by function here the
69 authors mean evidence of evolutionary selection. Not all functions are under evolutionary selection, consider
70 human olfactory receptors, many of which, although clearly functional, do accumulate deleterious mutations
71 and evolve almost neutrally. At the same time not all uORFs that exhibit $\omega \ll 1$ necessarily encode
72 functional peptides, because some uORFs are known to alter ribosome movement by making ribosomes stall
73 via specific interactions inside the peptide channel. Such stalling peptides may not function on their own outside
74 of the ribosome even though they would be expected to evolve as protein-coding.

75 **Response:** These comments are enlightening and much appreciated. To address these concerns, we performed
76 two additional analyses in this revised version. First, we performed PhyloCSF analysis of the coding regions
77 of uORFs. The PhyloCSF algorithm predicts whether a genomic region potentially represents a conserved
78 protein-coding region or not based on multiple sequence alignments³, and a positive PhyloCSF score means that
79 region is more likely to encode a peptide. As negative controls, we also calculated the PhyloCSF scores for
80 20,000 randomly selected ORFs in 3' UTRs (downstream ORFs, dORFs), as these dORFs have little chance of
81 translation. By comparing the PhyloCSF scores of the uORFs that showed evidence of translation with the
82 ribosome profiling data versus the random dORFs, we estimated that in humans, 0.44% (161 out of 36,655) of
83 the translated uORFs showed evidence of encoding conserved peptides, and that in *Drosophila*, 0.80% (102 of
84 12,754) of the translated uORFs might encode conserved peptides. Overall, these analyses suggest that less than
85 1% canonical uORFs might encode conserved peptides.

86 Next, we examined public mass spectrometry (MS) datasets for evidence of uORF-encoded peptides in
87 *Drosophila*. We analyzed the mass spectrometry (MS) data from 38 samples of different developmental stages
88 or tissues of *D. melanogaster* from previous studies⁴⁻⁸ (see Supplementary Table 7 for details). Among the
89 24,462 uORFs that met our parameter settings, 84 (0.34%) had peptides detected in at least one sample (see
90 Supplementary Table 8 for details). In combination with our finding that most uORFs do not encode conserved
91 peptides, these results suggest that only a very small fraction (< 1%) of the uORFs might encode peptides that
92 are maintained by natural selection during evolution.

93 The new results are described as follows (Page 10, Lines 318-338):

94 “To estimate the proportion of uORFs that might encode conserved peptides, for each uORF, we also
95 calculated PhyloCSF score, which predicts whether a genomic region potentially represents a conserved protein-
96 coding region or not based on multiple sequence alignments³ (a positive PhyloCSF score means that region is
97 more likely to encode a peptide). As a negative control, we also calculated the PhyloCSF scores for 20,000
98 randomly selected ORFs in 3' UTRs (downstream ORFs, dORFs), as these dORFs have little chance of
99 translation. Among the 36,655 uORFs that are ≥ 10 codons and evidenced of translation in humans, only 361
100 (0.985%) had positive PhyloCSF scores (Supplementary Fig. 12a). In contrast, the PhyloCSF score was positive
101 for 0.545% (109 out of 20,000) dORFs. Thus, after controlling for the background noises, only 0.44% (161) of
102 the translated uORFs showed evidence of encoding conserved peptides. In *Drosophila*, 1.19% (152 of 12,745)
103 translated uORFs and 0.39% (78 out of 20000 dORFs) had positive PhyloCSF scores, yielding an estimate of
104 0.80% (102 of 12,754) of the translated uORFs might encode conserved peptides. Overall, these analyses
105 suggest that less than 1% canonical uORFs might encode conserved peptides.

106 To test whether our evolutionary analyses of uORFs were supported by experimental evidence, we
107 analyzed the mass spectrometry (MS) data from 38 samples of different developmental stages or tissues of *D.*
108 *melanogaster* (Supplementary Table 8)⁴⁻⁸. Among the 24,462 uORFs that met our parameter settings (Methods),
109 84 (0.34%) had peptides detected in at least one sample (Supplementary Table 9). Interestingly, the BLS
110 analysis revealed that the MS-supported uORFs present more conserved coding regions than the other uORFs
111 (Fig. 5e), suggesting these MS-supported uORF peptides might be functionally important. Collectively, our
112 results support the notion that most uORFs play regulatory roles and their start codons are maintained due to

113 functional constraints, and only a tiny fraction (< 1%) of the uORFs might encode peptides that are maintained
114 by natural selection during evolution.”

115 We also put the following sentences in Discussion (Page 13, Lines 438-442), which are reproduced as follows:

116 “Overall, our results suggest that the major function of uORFs is to fine-tune CDS translation rather than
117 to encode conserved peptides. Nevertheless, we do not deny that some uORFs can encode functional peptides,
118 as clearly demonstrated by the previous studies⁹⁻¹¹. Of note, both our PhyloCSF analyses and MS data analyses
119 suggest that a small fraction (< 1%) of uORFs might produce peptides.”

120

121 Specific

122 1. It is not clear how exactly the groups of genes were divided into the categories for the analyses shown in Fig.
123 1Sb. A more detailed explicit description is necessary.

124 **Response:** We apologize that this information related to this figure (Supplementary Fig. 5b in the revised
125 manuscript) was not clearly described in our previous version. In this revised version, we presented the details
126 in the “Gene ontology analysis” subsection of Methods, which is reproduced as follows (Page 17, Lines 552-
127 560):

128 “Gene ontology annotations for human, mouse, rat, zebrafish, fly, *A. thaliana*, and yeast were downloaded
129 from the Gene Ontology Resource (2019-06-09 release). Because not all genes under a GO term were provided
130 in the GO annotation files, we parsed the gene annotation files to obtain the complete list of genes under each
131 term using topGO¹². For each species, all the GO terms belonging to Molecular Function (MF), Biological
132 Process (BP), and Cellular Component (CC) were combined in the enrichment analysis. The GO terms that
133 were enriched in uORF-containing genes or uORF-free genes were determined using Fisher's exact tests.
134 Multiple testing correction was performed with the Benjamini-Hochberg method¹³, and significant terms were
135 determined at a false discovery rate (FDR) of 0.1 for each species. Non-redundant representative terms that
136 were significantly enriched in at least five species were chosen for visualization.”

137 We also updated the figure legend (Supplementary Fig. 5b in the revised manuscript), which is reproduced
138 as follows:

139 “Gene categories enriched in uORF-free genes (left) and uORF-containing genes (right). In each species,
140 genes belonging to each category were extracted from the annotations provided by Gene Ontology Consortium
141 (Methods). Whether a gene category is enriched in the gene set is assessed with Fisher's exact test. Multiple
142 testing correction was performed with the Benjamini-Hochberg method¹³. The odds ratios (log2) and adjusted
143 *P* values are indicated by the color and size of the points, respectively. For uORF-containing genes, the same
144 analysis was performed for all the genes containing putative uORFs or only genes containing translated uORFs.
145 Non-redundant representative terms that are enriched in at least five of the seven model organisms were
146 displayed in the plot. See Supplementary Table 3 for the complete list of terms enriched in each species. Some
147 terms are insignificant in yeast, primarily because yeast is a unicellular organism with only 6,600 protein-coding
148 genes and 955 of those genes contain uORFs, which makes the statistical power of enrichment analyses
149 relatively low in yeast.”

150

151 2. The authors extensively used the data from ref. 27 (McGillivray et al) and attempted to make certain
152 conclusions regarding the evolution of uORFs reported in that work, for example, they found the evidence that
153 these uORFs are more conserved. This is inappropriate. McGillivray et al used conservation as one of the
154 features used for their machine learning algorithm: “Features were chosen to cover a broad range of categories
155 of data, including features associated with uORF position and length, conservation, functional metrics like RNA
156 expression, and sequence-based signatures that may relate to translation.” It makes no sense to show increased
157 conservation of uORFs that were predicted based on their conservation.

158 **Response:** Thank you for raising this concern. In our previous version, we mainly focused on the evolution of
159 canonical uORFs. We did not emphasize the evolutionary patterns of the noncanonical uORFs because 1) most
160 species surveyed in this study currently have no ribosome profiling data and, 2) it is very challenging to predict
161 the noncanonical uORFs *in silico* reliably. Hence, in our previous submission, we did not intend to show that
162 the noncanonical uORFs are more conserved. We touched the predicted noncanonical uORFs from McGillivray
163 *et al.*¹⁴ mainly in two places. First, when comparing the conservation of uAUGs relative to the background
164 triplets in the 5' UTRs, we excluded the start codons of 173,290 noncanonical uORFs (beginning with non-
165 AUG triplets) identified by McGillivray *et al.* from the backgrounds. Second, we found that the noncanonical
166 uORFs predicted in the previous study were indeed slightly more conserved than random triplets, but these are
167 significantly less conserved than the canonical uORFs.

168 We agree with the reviewer that it is inappropriate to compare the conservation of noncanonical uORFs
169 predicted based on conservation to the remaining random triplets. We rephrased the relevant sentences which
170 read as follows (Page 9, Lines 264-267): “We also calculated the BLS values for the start codons of the 173,290
171 noncanonical uORFs previously identified in humans by McGillivray *et al.*¹⁴. Since conservation was used as a
172 feature to identify the noncanonical uORFs in that study, it is not surprising that these noncanonical start codons
173 were slightly (~1.2 times) more conserved than the other random triplets ($P=2.1\times 10^{-77}$, WRST; Fig. 3d).
174 However, they were significantly less conserved than the canonical uAUGs ($P=1.0\times 10^{-10}$, WRST).”

175 Moreover, in the revised manuscript, we have added new analyses regarding the biological functions of
176 noncanonical uORFs from two aspects. First, we extracted previously published functional and population
177 genomic data and examined whether variations in uORF start codons influence the translation efficiency of the
178 main CDSs among different samples (Fig. 7a). Among the potentially functional uORFs in humans predicted
179 by McGillivray *et al.*¹⁴, 146 canonical and 796 noncanonical uORFs had genetic variants in their start codons
180 among these samples (only variants with minor allele frequency $\geq 5\%$ were considered in the analysis). We
181 performed linear regressions to assess the regulatory impact of uORF alteration on the translation of down-
182 stream CDSs, with a positive slope value in the regression meaning that the presence of a uORF in certain
183 individuals is associated with a decrease in the translation efficiency of the downstream CDS in those
184 individuals, and vice versa. A general trend was the slope values were overall positive for the canonical uORFs,
185 while the slope values for the noncanonical uORFs fluctuated around 0 (Fig. 6b). This comparison suggests that
186 in human populations, the noncanonical uORFs overall have relatively limited repressive effects on CDS
187 translation compared to the canonical uORFs.

188 Next, we experimentally verified the influence of both types of uORFs on CDS translation. We sampled

189 80 human uORFs and performed luciferase reporter assays in HEK293FT cells (Supplementary Fig. 17). These
190 tested uORFs, which included 42 canonical and 38 noncanonical ones, were predicted potentially functional by
191 McGillivray *et al.*¹⁴ and had polymorphic start codons in human populations. For each uORF, we compared the
192 repressive effect of the annotated uORF allele versus that of the non-uORF allele in suppressing the translation
193 of the reporter gene. Although occasionally the non-uORF allele had a stronger repressive effect than the uORF
194 allele, the general trend was that the uORF allele had a stronger effect than the non-uORF allele in suppressing
195 translation (Fig. 6c-d). Moreover, a significantly higher proportion of the canonical (55%, 23/42) than the
196 noncanonical (26%, 10/38) uORFs exhibited the pattern that the annotated uORF allele showed a significantly
197 stronger repressive effect on the CDS translation than the non-uORF allele ($P = 0.013$, Fisher's exact test, Fig.
198 6c-d). Also, the difference in CDS translation suppression between the uORF and the non-uORF allele is
199 significantly larger for the canonical than the noncanonical uORFs ($P = 0.006$, WRST). Altogether, these results
200 reinforced the thesis that the noncanonical uORFs overall have weaker repressive effects on CDS translation
201 than the canonical uORFs. We fully described these new results in the Section “**Comparing the canonical
202 versus noncanonical uORFs in repressing CDS translation in human populations**” (Page 12, Lines 379-
203 412).

204

205 3. Comparing the strength of Kozak context in uORFs and CDS ATGs. If we take two groups of sequences and
206 compare, they are likely to differ in some respect. If we define one as optimal, the other would become
207 suboptimal. Therefore, the purpose of this analysis is unclear to me. Perhaps it would be more meaningful to
208 compare three groups of contexts rather than two, by adding ATGs that are not used for initiation, e.g. internal
209 ATGs from CDS or ATGs from 3' UTRs or intergenic regions. We expect that the context of such ATGs should
210 not evolve to optimize translation initiation and would provide an estimate for a background context and a
211 variation in contexts. Then we would expect that uORFs context should be optimized for translation initiation,
212 but not as strong as CDS ATG. By having three points the authors could estimate whether uORF ATG context
213 is closer to neutral or that of CDS.

214 **Response:** Thank you for the constructive suggestions. In this revised version, we followed this reviewer's
215 suggestion and used the AUGs in 3' UTRs (dAUGs) as negative controls to discriminate whether the context of
216 uAUGs is optimized or close to the neutral background. We reported the new results as follows (Page 11, Lines
217 362-268):

218 “To test whether the sequence contexts of uAUGs are optimized, in each species we also calculated the
219 Kozak scores of the AUG triplets in 3' UTRs (downstream AUGs, dAUGs) as neutral controls. The Kozak
220 scores of uAUGs were significantly higher than those of dAUGs in most (88.2%, 120 out of 136) vertebrates,
221 (68.3%, 28 out of 41) plants and (180%, 180 out of 242) fungi; however, an opposite trend was observed in
222 invertebrates and no obvious trend was observed in protists (Supplementary Fig. 16b). These results suggest
223 that the optimization of the Kozak sequence context of uORFs is different across eukaryotic clades.”

224 We did not include the Kozak scores for the AUG triples inside the coding regions because the -6 to +4
225 nucleotides around such internal AUGs are under strong selective constraints due to coding functions or codon
226 usage bias.

227

228 4. “Unsurprisingly, for both uORFs and CDSs, the distance between two species from a clade tended to be
229 significantly shorter than that between one species in that clade and another species outside of that clade (Fig.
230 6c). These results suggest that the Kozak contextual characteristics tend to be similar between closely related
231 species for both uORFs and CDSs.”

232 This is indeed so unsurprising that it is unclear why was it even done. I believe that any other sequence, e.g. a
233 context of stop codons would exhibit the same behavior.

234 **Response:** We thank the reviewer for raising this concern. The aim of our analysis is to explore whether the
235 Kozak contexts around the start codons of uORFs are different across eukaryotic clades. To our knowledge,
236 such an issue has not been systematically explored yet. We think this question is important as there is growing
237 interest in engineering uORFs for precise translation control of the main protein products. Our results suggest
238 that considering species-specific Kozak sequence contextual features might be necessary in designing uORFs
239 for a specific desired trait in a certain species.

240 In this revised manuscript, we emphasized this point in Discussion (Page 14, Lines 468-471) with the
241 following sentences:” There has been a growing interest in engineering uORFs for precise translation control
242 of the main protein products¹⁵⁻¹⁷. Our results revealed the Kozak sequence context evolved across eukaryotic
243 clades, which suggests that the species-specific Kozak sequence contextual features should be considered in
244 designing uORFs for a specific desired trait.”

245

246

247 5. The authors made an observation that uORFs occurrence anticorrelates with expression levels. This makes
248 sense, but there could be at least two reasons for that. One is that the regulation usually works by suppression,
249 hence the mRNAs whose translation is regulated by uORFs are likely to be lowly expressed. The other is that
250 the negative selection acting on uAUGs is expected to be weaker for lowly expressed mRNAs. These two
251 scenarios are drastically different, could authors try to estimate contributions of each of these two scenarios?

252 **Response:** These comments are really insightful and enlightening. To address this issue, in this revised version,
253 we grouped genes of a species into 20 equal-sized bins based on increasing expression levels and calculated the
254 O/E ratio of uORFs in each bin. In all the five species we examined, the O/E ratio was substantially lower than
255 1 in each bin (Supplementary Fig. 5c), suggesting that purifying selection was the dominant evolutionary force
256 acting on the uORF occurrence regardless of gene expression levels. Nevertheless, we observed significant
257 anticorrelations between the gene expression level and O/E ratio of uORFs in each species, suggesting that
258 purifying selection acting on uAUGs is relatively weak for lowly expressed genes. On the other hand, genes in
259 certain functional categories, such as transcriptional factors, which are likely to be lowly expressed, might be
260 preferentially suppressed by uORFs at the translational level for optimizing protein production. Thus, our results
261 suggest that although gene expression level overall is an important factor influencing the genome-wide
262 distributions of uORFs across genes, the anticorrelation between gene expression level and uORF occurrence
263 is caused by very complex factors.

264 In this revised version, we have re-written the relevant section, which is reproduced as follows (Page 5,
265 Lines 139-167):

266 **“Gene expression level as an important factor influencing the genome-wide distributions of uORFs across**
267 **genes**

268 In humans, genes with uORFs exhibited lower expression levels than genes without uORFs¹⁸. Similarly, our
269 analysis of previously published mRNA and protein abundance data of fly, human, mouse, mustard plant, and
270 yeast revealed uORFs were infrequently detected in housekeeping genes, and there were significant
271 anticorrelations between the gene expression level and the number of uORFs (Supplementary Fig. 5a and
272 Supplementary Table 2). Meanwhile, gene ontology analysis revealed that genes containing putative uORFs
273 tend to be enriched in the categories of signal transduction, transcription factors, and membrane proteins
274 (Supplementary Fig. 5b; Supplementary Table 3). These patterns still held when we focused on the uORFs
275 supported by previously published ribosome profiling data in fly¹⁹ and other species collected in the GWIPs-
276 viz database²⁰ (Supplementary Table 4). Noteworthy, the anticorrelation between uORF occurrences and gene
277 expression level well reconciles with the gene ontology analyses as housekeeping genes tend to be highly (or
278 broadly) expressed²¹.

279 Since gene expression level affects the efficacy of natural selection²², we further asked whether the efficacy
280 of purifying selection is reduced in removing deleterious uORFs in lowly expressed genes. We grouped genes
281 of a species into 20 equal-sized bins based on increasing expression levels and calculated the O/E ratio of uORFs
282 in each bin. In all the five species we examined, the O/E ratio was lower than 1 in each bin (Supplementary Fig.
283 5c), suggesting that purifying selection was the dominant evolutionary force acting on the uORF occurrence
284 regardless of gene expression levels. Interestingly, we observed significant anticorrelations between the gene
285 expression level and O/E ratio of uORFs in each species, suggesting that purifying selection acting on uAUGs
286 is relatively weak for lowly expressed mRNAs.

287 Thus, our results suggest that gene expression level is an important factor influencing uORF distribution
288 across genes in a eukaryotic species. It is possible that excessive uORFs in highly expressed genes might cause
289 insufficient protein output, which is harmful to the organisms. We postulate that purifying selection has removed
290 deleterious uORFs in the highly expressed genes more efficiently than in the lowly expressed genes. On the
291 other hand, genes in certain functional categories, such as transcriptional factors, which are likely to be lowly
292 expressed, might be preferentially suppressed by uORFs at the translational level for optimizing protein
293 production. Further studies are needed to investigate the relative importance of the two mechanisms in shaping
294 the anticorrelation between gene expression level and uORF occurrence.”

295

296 6. To demonstrate the evidence of positive selection on 162 newly fixed uORFs, authors have used the
297 asymptotic McDonald-Kreitman, where the alpha parameter is the proportion of substitutions that are due to
298 adaptive evolution. But confidence intervals are quite wide and contain zeros (as well as negative values), so
299 there seems to be no strong evidence of positive selection (Fig 2a)

300 **Response:** We greatly appreciate this point. In this revision, we updated the analysis and made two interesting
301 observations. First, we detected strong and significant positive selection for newly fixed ATGs derived from

302 CpG to TpG mutations in primates. Second, we found that the signal of positive selection is more pronounced
303 in the genes with higher expression levels, both in primates and in *Drosophila*. In fact, the new results were
304 largely inspired by the enlightening comments this reviewer raised regarding the gene expression level in the
305 previous point. We reported the new results as follows (Page 7, Lines 210-225):

306 “We detected weak signals of positive selection on the newly fixed uORFs in all three branches, and the
307 value of α_{asym} , which represents the fraction of newly formed uORFs driven to fixation by positive selection,
308 was 0.18 (95% CI, -0.15~0.50), 0.15 (95% CI, -0.20~0.48), and 0.14 (95% CI, -0.21~0.48) in the three branches,
309 respectively (Fig. 2a). Noteworthy, C>T mutations at CpG dinucleotides are highly frequent in mammals²³, and
310 new AUGs can be generated from CpG to TpG mutations through two approaches²⁴: 1) from ACG to ATG, and
311 2) from CGTG to CATG (Fig. 2b). Thus, we further examined new uORFs derived from the CpG contexts and
312 the remaining new uORFs separately. Roughly speaking, ~33% of the new AUGs fixed in each of the three
313 branches were generated by CpG to TpG mutations. Interestingly, the CpG-derived uORFs were under strong
314 positive selection (the α_{asym} was 0.57 (95% CI, 0.29~0.85), 0.54 (95% CI, 0.24 ~ 0.83) and 0.53 (95% CI, 0.22~
315 0.83) in the three branches, respectively), while the α_{asym} for the remaining uORFs was close to 0 (Fig. 2b).
316 Noteworthy, the α_{asym} values were even higher when we focused on the new uORFs that were derived from the
317 CpG contexts in the highly expressed genes (Supplementary Table 6). Of note, for the new uORFs fixed in *D.*
318 *melanogaster* we previously analyzed¹⁹, a higher α_{asym} value was also observed for the highly expressed genes
319 (Supplementary Table 7). Therefore, although the prevalence of uORFs in a species was generally under
320 purifying selection, we still found a fraction of uORFs might be favored by positive selection even in primates
321 that typically have a small N_e .”

322

323 7. It is unclear how the relative fixation probability of newly originated uORFs was calculated. Could the authors
324 provide an explicit description of the procedure?

325 **Response:** We apologize for the obscure description of this process. In this revised manuscript, we have
326 provided a detailed description of the procedure used to calculate the relative fixation probability in “Methods”
327 (Page 14, Lines 464-472):

328 **“The fixation probability of new uORFs**

329 For a new autosomal mutation with a selective coefficient s in a diploid population of size N_e , the fixation
330 probability of the mutation relative to a neutral mutation was calculated as $f(s) =$

331 $2N_e \int_0^{\frac{1}{2N_e}} G(x)dx / \int_0^1 G(x)dx$, where $G(x) = \exp[-4N_e shx - 2N_e s(1 - 2h)x^2]$ and h is the dominance
332 coefficient²⁵. For mutations that introduce new uORFs into the population, the fractions of neutral, deleterious,
333 and beneficial mutations are denoted as p_1 , p_2 , and p_3 , respectively. Based on the assumption that the selective
334 coefficients for deleterious and beneficial mutations have the same absolute value, we can obtain the overall
335 relative fixation probability of mutations as $p_1 + p_2 f(-s) + p_3 f(s)$. In the simulation, we used a fixed $h=0.5$,
336 and p_1 , p_2 , and p_3 were set to 0.2, 0.75, and 0.05, respectively.”

337

338 **Reviewer #2 (Remarks to the Author):**

339 In the study of Zhang et al., the authors analyzed more than 10 million “uORFs” in over 200 eukaryotic species.
340 They found that 1) most of “uORFs” are under purifying selection. 2) the coding region of “uORFs” is overall
341 less conserved, suggesting that uORF is under neutral evolution or weak selective pressure. Finally, they also
342 analyzed the evolution of start codon and flanking context of uORFs. While the manuscript is written well,
343 many of main conclusions are not new, which have been reported by previous studies. Although previous studies
344 analyzed uORF evolution usually based on a small subset of closely-related species, simply using more species
345 does not significant extend our knowledge on the origin of uORF translation and its evolution. My major
346 concerns are as follows.

347 **Response:** We thank this reviewer for the thorough reviews. In this manuscript, we indeed found that 1) most
348 of “uORFs” are under purifying selection, and 2) the coding region of “uORFs” is overall less conserved, which
349 has been nicely summarized by this reviewer. **However, we went steps further than merely reporting such**
350 **observations.** The novel findings relevant to these two points are 1) we have demonstrated how positive and
351 purifying selection, coupled with differences in gene expression level and N_e , influence the genome-wide
352 distribution and contents of uORFs in eukaryotes, and 2) the uAUGs, particularly the translated uAUGs, tend
353 to be maintained by functional constraints during evolution, however, the coding regions of uORFs are overall
354 under neutral evolution. This comparison suggests that the major function of uORFs is to fine-tune CDS
355 translation rather than encode conserved peptides. In this new submission, inspired by the comments from this
356 reviewer, we also performed new analyses to highlight the novel findings of this study. For example, we further
357 carried out both phyloCSF analyses and mass spectrometry (MS) data analyses to demonstrate that only a small
358 fraction (< 1%) of uORFs might produce peptides (Page 10, Lines 318-338).

359 We believe the manuscript is much improved after addressing this reviewer’s concerns. We also highlighted
360 our changes in the Point-to-point response section.

361

362

363 1) Like canonical translation, uORF translation is energy-consuming. Uncontrolled uORF translation may
364 inhibit translation in main CDS. Therefore, it is not unexpected that the potential of uORF translation in 5’ UTR
365 has been eliminated during evolution. Also, similar observations have been reported by previous studies, for
366 example PMC5793785, PMC4890304.

367 **Response:** We agree with this reviewer that it is not surprising to find that uORFs are generally depleted in 5’
368 UTRs by natural selection. However, how the efficiency of purifying selection and positive selection acting on
369 uORFs prevalence differs across species or genes have not been thoroughly explored prior to this study.
370 Specifically, our novel discoveries are summarized as follows:

371 1) With the comparative analysis of uORF occurrences in 481 eukaryotes (242 fungi and 23 protists were newly
372 included in this revision), we found the trend of uORF depletion varies widely across different species, and the
373 degree of uORF depletion is mainly determined by the effective population size of a species.

374 2) Previously, we found that newly fixed uORFs in *D. melanogaster* are under positive selection. Here, we
375 extend this analysis to primates and found a similar trend, which suggests that although uORFs are overall under

376 depletion, a fraction of uORFs are favored by positive selection even in species with a small effective population
377 size such as primates. In particular, we detected strong and significant positive selection for newly fixed ATGs
378 derived from CpG to TpG mutations in primates. Moreover, we found that the signal of positive selection is
379 more pronounced in the genes with higher expression levels, both in primates and in *Drosophila*. We present
380 the updated results in Page 7, Lines 210-225.

381 3) We investigated the factors that influence the efficiency of uORF depletion within a species. We found that
382 the gene expression level is one major factor that influences uORF distribution across genes, and that purifying
383 selection on uORF prevalence is stronger in highly expressed genes. We present the new results in Page 6, Lines
384 152-167.

385 4) Based on the position and frame of a uORF relative to the downstream CDS, in this revision, we classified
386 uORFs into nonoverlapping uORFs, out-of-frame overlapping uORFs (oORFs), and N-terminal extensions. We
387 found that the O/E ratio of oORFs was overall significantly lower than that of nonoverlapping uORFs
388 (Supplementary Fig. 4). Interestingly, N-terminal extensions showed the lowest O/E ratio among the three
389 categories of uORFs in 460 out of 481 species (Supplementary Fig. 4), suggesting that novel N-terminal
390 extensions might be harmful to normal protein functions and tend to be depleted. We present the new results in
391 Page 5, Lines 122-128.

392 Given the above considerations, we believe that our study provides novel insights into the distribution and
393 evolution of uORFs in eukaryotes beyond the general idea that uORFs tend to be depleted from 5' UTR during
394 evolution.

395

396 2) uORF translation plays various roles in gene expression regulations. As demonstrated by many previous
397 studies (see reviews PMID: 28698598, PMID: 31003826), uORF may encode functional peptide, or uORF
398 translation may control downstream translation in main CDS. Again, it is not unexpected that the coding region
399 of uORFs may not under negative selection, if they do not encode functional peptides.

400 **Response:** Thank you for pointing this out. We agree with the reviewer that the coding region of an uORF is
401 likely not under negative selection if it does not encode functional peptides. However, the type of uORF
402 functions (regulatory versus coding) that dominate has not been thoroughly investigated. In this study, we
403 performed comparative studies of uORF start codons and coding regions through phylogenetic and population
404 genetic analyses. Our results suggest that although uORF start codons are more conserved than expected, the
405 coding regions of uORFs usually evolve neutrally or under weak selective constraints, which leads us to
406 conclude that most uORFs do not encode conserved peptides.

407 To estimate the proportion of uORFs that might encode conserved peptides, for each uORF, we also
408 calculated the PhyloCSF score, which predicts whether a genomic region potentially represents a conserved
409 protein-coding region or not based on multiple sequence alignments. Our PhyloCSF analysis revealed that only
410 0.44% of the human uORFs that are evidenced of translation might encode conserved peptides, and similarly,
411 0.80% (102 of 12,754) of the translated uORFs in *Drosophila* might encode conserved peptides. Overall, these
412 analyses suggest that less than 1% canonical uORFs might encode conserved peptides.

413 However, functional uORFs are not necessarily conserved. Therefore, we also searched for uORF-encoded
414 peptides among multiple mass spectrometry datasets in *D. melanogaster*. We found MS evidence for only 84
415 (0.23%) uORFs, and the peptides encoded by MS-supported uORFs are more conserved than those of the
416 remaining uORFs. Taken together, our results support the notion that the dominant function of uORFs is
417 regulatory rather than encoding peptides. The new analysis is described as follows (Page 10, Lines 318-338):

418 “To estimate the proportion of uORFs that might encode conserved peptides, for each uORF, we also
419 calculated PhyloCSF score, which predicts whether a genomic region potentially represents a conserved protein-
420 coding region or not based on multiple sequence alignments³ (a positive PhyloCSF score means that region is
421 more likely to encode a peptide). As a negative control, we also calculated the PhyloCSF scores for 20,000
422 randomly selected ORFs in 3' UTRs (downstream ORFs, dORFs), as these dORFs have little chance of
423 translation. Among the 36,655 uORFs that are ≥ 10 codons and evidenced of translation in humans, only 361
424 (0.985%) had positive PhyloCSF scores (Supplementary Fig. 12a). In contrast, the PhyloCSF score was positive
425 for 0.545% (109 out of 20,000) dORFs. Thus, after controlling for the background noises, only 0.44% (161) of
426 the translated uORFs showed evidence of encoding conserved peptides. In *Drosophila*, 1.19% (152 of 12,745)
427 translated uORFs and 0.39% (78 out of 20000 dORFs) had positive PhyloCSF scores, yielding an estimate of
428 0.80% (102 of 12,754) of the translated uORFs might encode conserved peptides. Overall, these analyses
429 suggest that less than 1% canonical uORFs might encode conserved peptides.

430 To test whether our evolutionary analyses of uORFs were supported by experimental evidence, we analyzed
431 the mass spectrometry (MS) data from 38 samples of different developmental stages or tissues of *D.*
432 *melanogaster* (Supplementary Table 8)⁴⁻⁸. Among the 24,462 uORFs that met our parameter settings (Methods),
433 84 (0.34%) had peptides detected in at least one sample (Supplementary Table 9). Interestingly, the BLS
434 analysis revealed that the MS-supported uORFs present more conserved coding regions than the other uORFs
435 (Fig. 5e), suggesting these MS-supported uORF peptides might be functionally important. Collectively, our
436 results support the notion that most uORFs play regulatory roles and their start codons are maintained due to
437 functional constraints, and only a tiny fraction ($< 1\%$) of the uORFs might encode peptides that are maintained
438 by natural selection during evolution.”

439 In Discussion (Page 13, Lines 438-442), we also revisit this point with the following sentences “Overall,
440 our results suggest that the major function of uORFs is to fine-tune CDS translation rather than to encode
441 conserved peptides. Nevertheless, we do not deny that some uORFs can encode functional peptides, as clearly
442 demonstrated by previous studies⁹⁻¹¹. Of note, both our PhyloCSF analyses and MS data analyses suggest that
443 a small fraction ($< 1\%$) of uORFs might produce peptides.”

444

445 3) About uORF definition. In this study, uORF is defined as a 5' UTR region starting with ATG and ending
446 with an in-frame stop codon (TAG, TAA or TGA). The uORF definition is problematic. First, they overlooked
447 uORFs starting with non-canonical start codons such as CTG, TTG or ATT. Previous studies have suggested
448 that non-canonical start codons are more prevalence than canonical start codon (i.e. ATG) in uORFs. Second, I
449 would define these regions as putative uORFs or potential uORFs, because majority of these so-called “uORFs”
450 are not translatable. Only a very small number of these putative uORFs are real uORFs with significant protein

451 translation. Analysis based on these putative uORFs will be strongly affected by huge amount of false positives
452 (or background noises), and can not be used to support the conclusions on uORFs.

453 **Response:** Thank you for raising this concern. We apologize for not clearly stating that the focus of this study
454 was canonical uORFs that start with AUG. As most species surveyed in this study currently have no ribosome
455 profiling data, and it is very challenging to predict the noncanonical uORFs *in silico* reliably, we only focused
456 on the putative canonical uORFs which start with the AUG start codon. In this revised manuscript, we have
457 inserted the following sentences (Lines 73-77): “As most species surveyed in this study currently have no
458 ribosome profiling data, and it is very challenging to predict the noncanonical uORFs *in silico* reliably, we only
459 focused on the putative canonical uORFs which start with the AUG start codon. Hence, in what follows, the
460 uORFs analyzed in this study are restricted to the putative canonical uORFs unless explicitly stated otherwise
461 (all the annotated uORFs are presented in figshare²⁶).”

462 For the putative canonical uORFs, recent studies suggest that most of them showed evidence of translation,
463 although the signal of translation is dependent on the sequencing coverage of Ribo-Seq and the number of
464 samples surveyed in a species. For example, based on the currently available ribosome profiling data from
465 humans, mice, and flies, we found that approximately 70-90% of canonical uORFs can be translated
466 (Supplementary Table 3). Moreover, our previous analysis suggested that many uORFs are selectively used
467 during development or in different tissues¹⁹, which suggested that more translated uORFs can be found if we
468 profile more developmental stages, tissues, or cell lines from a species. Therefore, by focusing on the putative
469 canonical uORFs only, we can limit the influence of potential false positives.

470 We agree with the reviewer that only a small number of noncanonical uORFs might be real, considering
471 the large number of putative non-AUG uORFs in the genomes (approximately 1.2 million in humans, as reported
472 by McGillivray et al.). In the revised manuscript, we have added new analyses regarding the biological functions
473 of noncanonical uORFs from two aspects. First, we extracted previously published functional and population
474 genomic data and examined whether variations in uORF start codons influence the translation efficiency of the
475 main CDSs among different samples (Fig. 7a). Among the potentially functional uORFs in humans predicted
476 by McGillivray *et al.*¹⁴, 146 canonical and 796 noncanonical uORFs had genetic variants in their start codons
477 among these samples (only variants with minor allele frequency $\geq 5\%$ were considered in the analysis). We
478 performed linear regressions to assess the regulatory impact of uORF alteration on the translation of down-
479 stream CDSs, with a positive slope value in the regression meaning that the presence of a uORF in certain
480 individuals is associated with a decrease in the translation efficiency of the downstream CDS in those
481 individuals, and vice versa. A general trend was the slope values were overall positive for the canonical uORFs,
482 while the slope values for the noncanonical uORFs fluctuated around 0 (Fig. 6b). This comparison suggests that
483 in human populations, the noncanonical uORFs overall have relatively limited repressive effects on CDS
484 translation compared to the canonical uORFs.

485 Next, we experimentally verified the influence of both types of uORFs on CDS translation, we sampled 80
486 human uORFs and performed luciferase reporter assays in HEK293FT cells (Supplementary Fig. 17). These
487 tested uORFs, which included 42 canonical and 38 noncanonical ones, were predicted potentially functional by
488 McGillivray *et al.*¹⁴ and had polymorphic start codons in human populations. For each uORF, we compared the

489 repressive effect of the annotated uORF allele versus that of the non-uORF allele in suppressing the translation
490 of the reporter gene. Although occasionally the non-uORF allele had a stronger repressive effect than the uORF
491 allele, the general trend was that the uORF allele had a stronger effect than the non-uORF allele in suppressing
492 translation (Fig. 6c-d). Moreover, a significantly higher proportion of the canonical (55%, 23/42) than the
493 noncanonical (26%, 10/38) uORFs exhibited the pattern that the annotated uORF allele showed a significantly
494 stronger repressive effect on the CDS translation than the no-uORF allele ($P = 0.013$, Fisher's exact test, Fig.
495 6c-d). Also, the difference in CDS translation suppression between the uORF and the non-uORF allele is
496 significantly larger for the canonical than the noncanonical uORFs ($P = 0.006$, WRST). Altogether, these results
497 reinforced the thesis that the noncanonical uORFs overall have weaker repressive effects on CDS translation
498 than the canonical uORFs. We fully described these new results in the Section “**Comparing the canonical
499 versus noncanonical uORFs in repressing CDS translation in human populations**” (Page 14, Lines 384-
500 412). The new results are described as follows:

501 “To test whether the noncanonical uORFs influence the translation of CDSs, we extracted high-quality
502 genotyping, mRNA-Seq, and Ribo-Seq data of 60 human lymphoblastoid cell lines from previous studies^{27,28},
503 and examined whether variations in uORF start codons influence the translation efficiency of the main CDSs
504 among different samples (Fig. 7a). Among the potentially functional uORFs in humans predicted by
505 McGillivray *et al.*¹⁴, 146 canonical and 796 noncanonical uORFs had genetic variants in their start codons
506 among these samples (only variants with minor allele frequency $\geq 5\%$ were considered in the analysis). We
507 performed linear regressions to assess the regulatory impact of uORF alteration on the translation of down-
508 stream CDSs, with a positive slope value in the regression meaning that the presence of a uORF in certain
509 individuals is associated with a decrease in the translation efficiency of the downstream CDS in those
510 individuals, and vice versa (Methods). A general trend was the slope values were overall positive for the
511 canonical uORFs, while the slope values for the noncanonical uORFs fluctuated around 0 (Fig. 6b). This
512 comparison suggests that in human populations, the noncanonical uORFs overall have relatively limited
513 repressive effects on CDS translation compared to the canonical uORFs, although we cannot exclude the
514 possibility that a small fraction of the noncanonical uORFs might have strong repressive effects on the
515 translation of downstream CDSs.

516 To experimentally verify the influence of both types of uORFs on CDS translation, we sampled 80 human
517 uORFs and performed luciferase reporter assays in HEK293FT cells (Supplementary Fig. 17). These tested
518 uORFs, which included 42 canonical and 38 noncanonical ones, were predicted potentially functional by
519 McGillivray *et al.*¹⁴ and had polymorphic start codons in human populations. For each uORF, we compared the
520 repressive effect of the annotated uORF allele versus that of the non-uORF allele in suppressing translation of
521 the reporter gene. Although occasionally the non-uORF allele had a stronger repressive effect than the uORF
522 allele, the general trend was that the uORF allele had a stronger effect than the non-uORF allele in suppressing
523 translation (Fig. 6c-d). Moreover, a significantly higher proportion of the canonical (55%, 23/42) than the
524 noncanonical (26%, 10/38) uORFs exhibited the pattern that the annotated uORF allele showed a significantly
525 stronger repressive effect on the CDS translation than the non-uORF allele ($P = 0.013$, Fisher's exact test, Fig.
526 6c-d). Also, the difference in CDS translation suppression between the uORF and the non-uORF allele is

527 significantly larger for the canonical than the noncanonical uORFs ($P = 0.006$, WRST). Altogether, these results
528 reinforced the thesis that the noncanonical uORFs overall have weaker repressive effects on CDS translation
529 than the canonical uORFs.”

530

531 For example,

532 i) they found O/E ratio (based on these putative uORFs) is significant lower than 1, suggesting that “purifying
533 selection is the major force shaping the prevalence of uORFs”. This result only suggests that ATG triplets are
534 depleted in 5’ UTR. Purifying selection for ATG triplets in 5’ UTR does not mean a necessary of selection for
535 uORFs. In fact, at least in yeast, a previous study (PMC5793785) reported an elevated non-canonical start codon
536 in 5’ UTR, indicating a possibility to maintain some kinds of uORF translation.

537 Response: Thank you for pointing this out. We apologize that in our previous submission, we did not clearly
538 explain how natural selection, i.e., positive and purifying selection, coupled with differences in gene expression
539 level and N_e , influence the genome-wide distribution and contents of uORFs in eukaryotes. As explained above,
540 in our analysis, we mainly focused on the putative canonical uORFs. In 475 out of 481 species we analyzed
541 (216 multicellular plants and animals, and 242 fungi and 23 protists that were added in the revised manuscript),
542 the O/E ratio of putative uORFs was significantly lower than 1, suggesting the prevalence of canonical uORFs
543 in a species was generally under purifying selection. However, we also found that positive selection can drive
544 the fixation of new uORFs that are beneficial in primates that typically have a small N_e . These results suggest
545 that positive selection might play a more important role in driving uORF evolution than previously anticipated.
546 Furthermore, we demonstrated that how the effective population size of a species affects the efficacy of natural
547 selection on the prevalence of uORFs. The whole section is presented in Lines 202-239 of Pages 7-8.

548 Moreover, in this revised manuscript, we further explored how gene expression level is an important factor
549 influencing the distribution of uORFs across genes. Our gene ontology analysis revealed that uORFs are biased
550 in genes of different functional categories, which are associated with gene expression levels. We also found that
551 purifying selection has removed deleterious uORFs in the highly expressed genes more efficiently than in the
552 lowly expressed genes. We presented the new analysis in the section “**Gene expression level as an important
553 factor influencing the genome-wide distributions of uORFs across genes**” (Page 5, Lines 139-167). We also
554 showed that the efficacy of positive selection on uORFs is stronger in highly expressed genes, and this pattern
555 was observed in both primates and *Drosophila* (Page 7, Lines 220-223).

556 Also, in Discussion, we reconciled the roles of natural selection and gene expression level on the
557 distribution of uORFs across species and across genes. The relevant sentences are reproduced as follows (Page
558 13, Lines 416-431):

559 “Although the prevalence of canonical uORFs in a species was generally under purifying selection, we
560 still found a fraction of new canonical uORFs might be favored by positive selection even in primates that
561 typically have a small N_e . These observations are consistent with the evolution model of uORFs we previously
562 proposed^{19,29}. Under that model, the majority of newly formed uORFs are deleterious and quickly removed from
563 the population, and a relatively smaller fraction of the new uORFs are beneficial and rapidly fixed in populations
564 under positive selection. After fixation, the functional uORFs, particularly the start codons, are maintained by

565 natural selection during evolution. Hence, although in a species the occurrence of a uORF is influenced by
566 positive or purifying selection, the opposing effects of positive selection and purifying selection acting on new
567 uORFs result in a pattern that uORFs are overall depleted in 5' UTRs. As shown in our population genetic
568 modeling, the efficacies of both positive and purifying selection on uORF fixation in a species are influenced
569 by the effective population size. Moreover, we also found that gene expression level affects the efficiency of
570 natural selection acting on uORF occurrences. Thus, our results have systematically demonstrated how positive
571 and purifying selection, coupled with differences in gene expression level and N_e , influence the genome-wide
572 distribution and contents of uORFs in eukaryotes. Together, our analyses provide an unprecedented overview
573 of the general principles underlying the distribution and sequence evolution of uORFs in eukaryotes.”

574

575 ii) The authors found that the dN/dS ratio for uORF CDS is “roughly equal to 1 between human and macaque”.
576 They concluded that this result supports neutral evolution of uORFs. However, because majority of “uORFs”
577 in their datasets are non-translatable (or not real uORF), these negative uORFs may significantly increase the
578 dN/dS ratio, since they encode nothing. Again, in *Drosophila*, the dN/dS ratio for all uORF CDSs is close to 1,
579 but later, they found that “uORFs with higher Kozak scores presented significantly lower dN/dS ratio in
580 *Drosophila*, suggesting a scenario in which the coding regions of uORFs with optimal Kozak sequence context
581 are under stronger purifying selection in *Drosophila*”. Since ATG surrounded by Kozak sequences are more
582 likely to be translated, I believe their negative result (i.e. dN/dS is close to 1) is due to too many negative uORFs
583 in their datasets.

584 **Response:** Thank you for raising these concerns. We apologize that in our previous submission, we might have
585 performed the analyses in an unnecessarily complicated approach so that our results might have been misleading
586 to this reviewer. In this study, we mainly focused on the putative canonical uORFs, and 69% and 89% of such
587 uORFs exhibit evidence of translation in humans and flies, respectively. Therefore, our results are not likely
588 caused by too many negative uORFs as this reviewer thought. Furthermore, we also repeated all the analyses in
589 this section using the uORFs that showed evidence of translation (Supplementary Fig. 11), and our conclusions
590 were not affected. In this new submission, we updated our analyses to avoid potential misunderstandings.

591 Moreover, in this revised version, we performed two additional analyses to address this reviewer’s
592 concerns. First, we performed PhyloCSF analysis of the coding regions of uORFs. The PhyloCSF algorithm
593 predicts whether a genomic region potentially represents a conserved protein-coding region or not based on
594 multiple sequence alignments³, and a positive PhyloCSF score means that region is more likely to encode a
595 peptide. As negative controls, we also calculated the PhyloCSF scores for 20,000 randomly selected ORFs in 3’
596 UTRs (downstream ORFs, dORFs), as these dORFs have little chance of translation. By comparing the
597 PhyloCSF scores of the uORFs that showed evidence of translation with the ribosome profiling data versus the
598 random dORFs, we estimated that in humans, 0.44% (161 out of 36,655) of the translated uORFs showed
599 evidence of encoding conserved peptides, and that in *Drosophila*, 0.80% (102 of 12,754) of the translated
600 uORFs might encode conserved peptides. Overall, these analyses suggest that less than 1% canonical uORFs
601 might encode conserved peptides.

602 Next, we examined public mass spectrometry (MS) datasets for evidence of uORF-encoded peptides in
603 *Drosophila*. We analyzed the mass spectrometry (MS) data from 38 samples of different developmental stages
604 or tissues of *D. melanogaster* from previous studies⁴⁻⁸ (see Supplementary Table 8 for details). Among the
605 24,462 uORFs that met our parameter settings, 84 (0.34%) had peptides detected in at least one sample (see
606 Supplementary Table 9 for details). In combination with our finding that most uORFs do not encode conserved
607 peptides, these results suggest that only a very small fraction ($< 1\%$) of the uORFs might encode peptides that
608 are maintained by natural selection during evolution. Therefore, we obtained consistent results among our
609 molecular evolution and population genetic analysis, the phyloCSF analysis, and the MS data re-analysis. The
610 new results are described in Lines 318-338 of Page 10.

611
612 iii) the same problem can be found in the analysis of “evolution of contextual characteristics that influence
613 uORF translation”.

614 **Response:** The effect of potential negative uORFs should be limited since we only focused on canonical uORFs,
615 most of which are evidenced of translation with the ribosome profiling data in model organisms.

616
617 4) Page 4, line 94. “gene expression level is a major determinant of the uORF distribution across genes in a
618 eukaryotic species” Because the number of putative uORFs positively correlates with 5' UTR length, I
619 wondered whether 5' UTR may confound the correlation of putative uORF number to gene expression.

620 **Response:** Thank you for this insightful comment. To address this concern, in this revised version, we grouped
621 genes of a species into bins of equal size based on their expression level and calculated the O/E ratio of uORFs
622 for genes in each bin. As a result, the potential confounding effect of differences in the 5' UTR length is also
623 properly controlled. In all the five species we examined, the O/E ratio was lower than 1 in each bin
624 (Supplementary Fig. 5c), suggesting that purifying selection was the dominant evolutionary force acting on the
625 uORF occurrence regardless of gene expression levels. Interestingly, we observed significant anticorrelations
626 between the gene expression level and O/E ratio of uORFs in each species, suggesting that purifying selection
627 acting on uAUGs is relatively weak for lowly expressed mRNAs. The new analysis is described as follows
628 (Page 6, Lines 152-167):

629 “Since gene expression level affects the efficacy of natural selection²², we further asked whether the
630 efficacy of purifying selection is reduced in removing deleterious uORFs in lowly expressed genes. We grouped
631 genes of a species into 20 equal-sized bins based on increasing expression levels and calculated the O/E ratio
632 of uORFs in each bin. In all the five species we examined, the O/E ratio was lower than 1 in each bin
633 (Supplementary Fig. 5c), suggesting that purifying selection was the dominant evolutionary force acting on the
634 uORF occurrence regardless of gene expression levels. Interestingly, we observed significant anticorrelations
635 between the gene expression level and O/E ratio of uORFs in each species, suggesting that purifying selection
636 acting on uAUGs is relatively weak for lowly expressed mRNAs.

637 Thus, our results suggest that gene expression level is an important factor influencing uORF distribution
638 across genes in a eukaryotic species. It is possible that excessive uORFs in highly expressed genes might cause
639 insufficient protein output, which is harmful to the organisms. We postulate that purifying selection has removed

640 deleterious uORFs in the highly expressed genes more efficiently than in the lowly expressed genes. On the
641 other hand, genes in certain functional categories, such as transcriptional factors, which are likely to be lowly
642 expressed, might be preferentially suppressed by uORFs at the translational level for optimizing protein
643 production. Further studies are needed to investigate the relative importance of the two mechanisms in shaping
644 the anticorrelation between gene expression level and uORF occurrence.”

645

646 5) Page 4, line 101, “while maintaining the same dinucleotide frequency”. Please explain why dinucleotide
647 frequency is maintained. Does single, or trip-nucleotide frequency significantly affect O/E ratio?

648 **Response:** Thank you for the helpful suggestion. In this revision, we further explained the reason for
649 maintaining the same dinucleotide frequency as follows (Page 4, Lines 107-110): “We maintained the same
650 dinucleotide frequencies in each sequence during shuffling for two reasons. First, the stacking energy of a new
651 base pair is influenced by the neighboring base pairs in an RNA molecule^{30,31}. Second, the biased mutations in
652 certain dinucleotide contexts, such as from CpG to TpG mutations in mammals, might also affect the prevalence
653 of uORFs.”

654 Since the AUG start codon has three nucleotides, maintaining the trinucleotide frequency means that the
655 frequency of every triplet in a sequence is unchanged in the shuffled sequence. Therefore, the O/E ratio for any
656 triplet (including ATG) will be 1, which is not appropriate for the current study.

657

658 6) O/E ratio in 5' UTR might be ok to estimate the selection for ATG triplets. To strength the results, O/E ratio
659 in 3' UTR should be considered as negative control, since translation in 3' UTR ORFs is less likely than that in
660 5' UTR. In addition, it would be great if O/E ratios for the other 61 triplets are displayed.

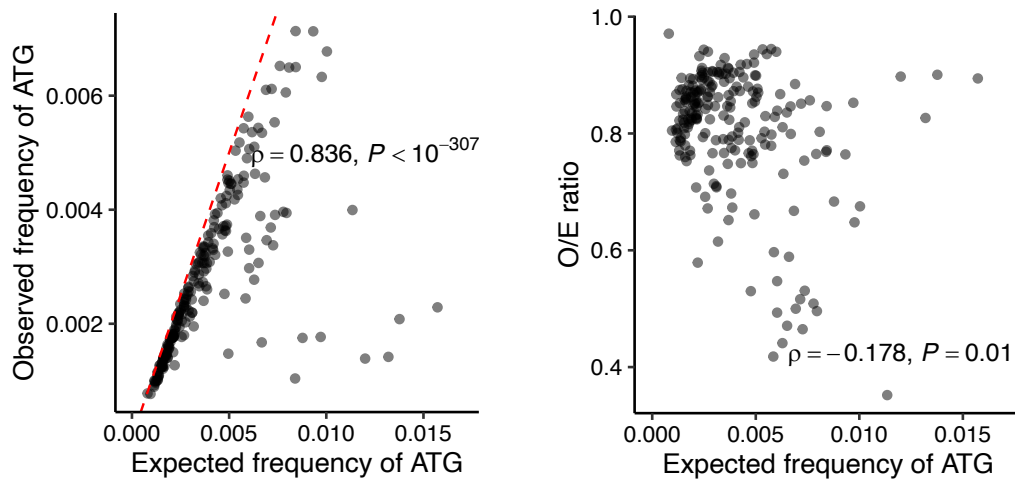
661 **Response:** Thank you for these helpful suggestions. In this revision, we followed these suggestions and
662 calculated the O/E ratios for all triplets in the 5' UTR and 3' UTR. The new results have been included in the
663 revised manuscript, and relevant sentences read as follows (Page 4, Lines 113-117):

664 “Since AUG is the defining feature of a canonical uORF, the O/E ratio is essentially the observed/expected
665 number of AUG triplets in the 5' UTRs. As a negative control, we also calculated the O/E ratio of all the other
666 63 possible triplets in 5' UTRs and 3' UTRs separately in each species. Of note, AUG had the lowest relative
667 O/E ratio (5' UTRs over 3' UTRs) among all the 64 possible triplets (Supplementary Fig. 2), supporting the
668 notion that purifying selection is the major force shaping the prevalence of uORFs in the eukaryotic genomes.”

669

670 7) Page 5, line 122. “The O/E ratio varied wildly across the 216 species”, is this ratio affected by different
671 background ATG frequency (E) across the species?

672 **Response:** Thank you for asking this question. Overall, the background ATG frequency (E) has little influence
673 on the O/E ratio of ATG in the 5' UTRs among different species. E is highly correlated with the observed
674 frequency of ATG (O) in a species (Fig. R1 left). In contrast, the O/E ratio showed a much weaker correlation
675 with the background ATG frequency (E), which is likely due to the uneven distribution of background ATG
676 frequencies across species. Moreover, the O/E ratio enabled the efficient measurement of selective pressure on
677 uORF depletion in a given species, as shown in previous studies^{19,32-35}.



678

679 Fig. R1. Influence of the expected frequency of ATG. Left panel, the relationship between the expected
 680 frequency of ATG and the O/E ratio among different species. Right panel, the relationship between the expected
 681 frequency of ATG and the observed frequency of ATG. Spearman's correlation analysis was performed for
 682 each plot.

683

684 8) Page 8, they found longer uORFs have fewer conserved peptides. This is a little unexpected to me. Because
 685 uORF translation is energy-consuming. If a uORF plays regulator role, a shorter ORF is sufficient to block
 686 ribosome scanning to downstream region. The longer ORF does not significant benefit the regulator role, but
 687 indeed consume more energy.

688 **Response:** Thank you for the suggestion. We have revised the relevant text to reflect this point (Page 10, Lines
 689 291-296):

690 "Of note, a strong anticorrelation was observed between the BLSs and the lengths of uORF peptides in
 691 both humans and flies (see Fig. 4c and 4d), suggesting the peptides encoded by long uORFs are less likely to
 692 be maintained during evolution because they were more likely disrupted by stop codons or frameshifts. Also, if
 693 the major function of uORFs is to regulate CDS translation, a longer uORF might be less advantageous than a
 694 shorter one because the translation of a longer uORF consumes more energy and metabolites, which might be
 695 harmful to the host organisms."

696

697

698

699 **References:**

- 700 1 Sengupta, S. & Higgs, P. G. Pathways of Genetic Code Evolution in Ancient and Modern Organisms. *J*
 701 *Mol Evol* **80**, 229-243, doi:10.1007/s00239-015-9686-8 (2015).
 702 2 Swart, E. C., Serra, V., Petroni, G. & Nowacki, M. Genetic Codes with No Dedicated Stop Codon:
 703 Context-Dependent Translation Termination. *Cell* **166**, 691-702, doi:10.1016/j.cell.2016.06.020
 704 (2016).
 705 3 Lin, M. F., Jungreis, I. & Kellis, M. PhyloCSF: a comparative genomics method to distinguish protein
 706 coding and non-coding regions. *Bioinformatics* **27**, i275-282, doi:10.1093/bioinformatics/btr209
 707 (2011).
 708 4 Xing, X. *et al.* Qualitative and quantitative analysis of the adult *Drosophila melanogaster* proteome.
 709 *Proteomics* **14**, 286-290, doi:10.1002/pmic.201300121 (2014).
 710 5 Casas-Vila, N. *et al.* The developmental proteome of *Drosophila melanogaster*. *Genome Res* **27**, 1273-

1285, doi:10.1101/gr.213694.116 (2017).

712 6 Ashley, J. *et al.* Retrovirus-like Gag Protein Arc1 Binds RNA and Traffics across Synaptic Boutons. *Cell* **172**, 262-274 e211, doi:10.1016/j.cell.2017.12.022 (2018).

713 7 Kuznetsova, K. G. *et al.* Proteogenomics of Adenosine-to-Inosine RNA Editing in the Fruit Fly. *J Proteome Res* **17**, 3889-3903, doi:10.1021/acs.jproteome.8b00553 (2018).

714 8 Sabbadin, F. *et al.* An ancient family of lytic polysaccharide monooxygenases with roles in arthropod development and biomass digestion. *Nat Commun* **9**, 756, doi:10.1038/s41467-018-03142-x (2018).

715 9 Mackowiak, S. D. *et al.* Extensive identification and analysis of conserved small ORFs in animals. *Genome Biol* **16**, 179, doi:10.1186/s13059-015-0742-x (2015).

716 10 Aspden, J. L. *et al.* Extensive translation of small Open Reading Frames revealed by Poly-Ribo-Seq. *eLife* **3**, e03528, doi:10.7554/eLife.03528 (2014).

717 11 van der Horst, S., Snel, B., Hanson, J. & Smeekens, S. Novel pipeline identifies new upstream ORFs and non-AUG initiating main ORFs with conserved amino acid sequences in the 5' leader of mRNAs in *Arabidopsis thaliana*. *Rna* **25**, 292-304, doi:10.1261/rna.067983.118 (2019).

718 12 Alexa, A. & Rahnenfuhrer, J. topGO: enrichment analysis for gene ontology. *R package* (2019).

719 13 Benjamini, Y. & Hochberg, Y. Controlling the false discovery rate: a practical and powerful approach to multiple testing. *Journal of the Royal statistical society: series B (Methodological)* **57**, 289-300 (1995).

720 14 McGillivray, P. *et al.* A comprehensive catalog of predicted functional upstream open reading frames in humans. *Nucleic acids research* **46**, 3326-3338, doi:10.1093/nar/gky188 (2018).

721 15 Xu, G. *et al.* uORF-mediated translation allows engineered plant disease resistance without fitness costs. *Nature* **545**, 491-494, doi:10.1038/nature22372 (2017).

722 16 Zhang, H. *et al.* Genome editing of upstream open reading frames enables translational control in plants. *Nature Biotechnology* **36**, 894-898, doi:10.1038/nbt.4202 (2018).

723 17 Ferreira, J. P., Overton, K. W. & Wang, C. L. Tuning gene expression with synthetic upstream open reading frames. *Proceedings of the National Academy of Sciences* **110**, 11284, doi:10.1073/pnas.1305590110 (2013).

724 18 Ye, Y. *et al.* Analysis of human upstream open reading frames and impact on gene expression. *Hum Genet* **134**, 605-612, doi:10.1007/s00439-015-1544-7 (2015).

725 19 Zhang, H. *et al.* Genome-wide maps of ribosomal occupancy provide insights into adaptive evolution and regulatory roles of uORFs during *Drosophila* development. *PLoS Biol* **16**, e2003903, doi:10.1371/journal.pbio.2003903 (2018).

726 20 Michel, A. M., Kiniry, S. J., O'Connor, Patrick B F., Mullan, J. P. & Baranov, P. V. GWIPS-viz: 2018 update. *Nucleic Acids Research* **46**, D823-D830, doi:10.1093/nar/gkx790 (2017).

727 21 Eisenberg, E. & Levanon, E. Y. Human housekeeping genes, revisited. *Trends in Genetics* **29**, 569-574, doi:10.1016/j.tig.2013.05.010 (2013).

728 22 Zhang, J. & Yang, J.-R. Determinants of the rate of protein sequence evolution. *Nature reviews. Genetics* **16**, 409-420, doi:10.1038/nrg3950 (2015).

729 23 Gonzalez-Perez, A., Sabarinathan, R. & Lopez-Bigas, N. Local Determinants of the Mutational Landscape of the Human Genome. *Cell* **177**, 101-114, doi:10.1016/j.cell.2019.02.051 (2019).

730 24 Kitano, S., Kurasawa, H. & Aizawa, Y. Transposable elements shape the human proteome landscape via formation of cis-acting upstream open reading frames. *Genes Cells* **23**, 274-284, doi:10.1111/gtc.12567 (2018).

731 25 Kimura, M. Diffusion models in population genetics. *Journal of Applied Probability* (1964).

732 26 Hong, Z. *et al.* The annotation of uORFs in 481 eukaryotes. *figshare*, doi:10.6084/m9.figshare.9980441.v2 (2020).

733 27 Auton, A. *et al.* A global reference for human genetic variation. *Nature* **526**, 68-74, doi:10.1038/nature15393 (2015).

734 28 Battle, A. *et al.* Genomic variation. Impact of regulatory variation from RNA to protein. *Science* **347**, 664-667, doi:10.1126/science.1260793 (2015).

735 29 Zhang, H., Wang, Y. & Lu, J. Function and Evolution of Upstream ORFs in Eukaryotes. *Trends in Biochemical Sciences* **44**, 782-794, doi:10.1016/j.tibs.2019.03.002 (2019).

736 30 Clote, P., Ferré, F., Kranakis, E. & Krizanc, D. Structural RNA has lower folding energy than random RNA of the same dinucleotide frequency. *RNA (New York, N.Y.)* **11**, 578-591, doi:10.1261/rna.7220505 (2005).

737 31 Workman, C. & Krogh, A. No evidence that mRNAs have lower folding free energies than random sequences with the same dinucleotide distribution. *Nucleic acids research* **27**, 4816-4822,

768 doi:10.1093/nar/27.24.4816 (1999).
769 32 Kozak, M. Possible role of flanking nucleotides in recognition of the AUG initiator codon by eukaryotic
770 ribosomes. *Nucleic Acids Res* **9**, 5233-5252 (1981).
771 33 Lynch, M., Scofield, D. G. & Hong, X. The evolution of transcription-initiation sites. *Mol Biol Evol* **22**,
772 1137-1146, doi:10.1093/molbev/msi100 (2005).
773 34 Neafsey, D. E. & Galagan, J. E. Dual Modes of Natural Selection on Upstream Open Reading Frames.
774 *Molecular Biology and Evolution* **24**, 1744-1751, doi:10.1093/molbev/msm093 (2007).
775 35 Rogozin, I. B., Kochetov, A. V., Kondrashov, F. A., Koonin, E. V. & Milanesi, L. Presence of ATG
776 triplets in 5' untranslated regions of eukaryotic cDNAs correlates with a 'weak' context of the start
777 codon. *Bioinformatics* **17**, 890-900 (2001).

Reviewers' Comments:

Reviewer #1:

Remarks to the Author:

The authors performed several analyses which expanded the manuscript substantially and well beyond my expectations. Because of such an extensive revision, a few issues emerged that are specific for the revised version, which I mention below. However, the significance of these issues is relatively minor. I think that the general message of this manuscript is adequately supported.

Comments pertinent for the revised version.

My very first comment related to the apparent impossibility of uORF (if we define uORFs as short translated sequences upstream of CDS) existence in some protists where genetic codes do not allow translation termination of translation in the internal positions of mRNA.

In these organisms any translation initiation event would result in a production of a long protein. So I suggested that the authors should simply acknowledge this fact and modify their title and discussion to make it more specific, i.e. plants and animals, but not eukaryotes in general.

Instead of following this simple suggestion, the authors chose a hard way, they decided to expand their analysis to include fungi and protists. Of course, this broadens the manuscript and I agree that the use of eukaryotes in the title is appropriate now. However, it made the work not only broader, but also more complicated. While the situation with fungi is similar to that of plants and animals, protists seem to be different and I would like to ask authors to make additional changes in their manuscript to clarify the situation.

First, protists are hugely diverse phylogenetically and exhibit considerable diversity in organisation of their genetic information in their genomes and probably in mRNAs as well. Because of that the analysis of 23 species may not be adequate to obtain a general picture.

Second, the analysis was done incorrectly. The authors did not take into account the diversity of their genetic codes in these species: "Putative uORFs that start with AUG codons and end with stop codons (UAA/UAG/UGA) were identified from the annotated 5' UTRs of protein-coding genes". This is clearly wrong for the organisms that do not use one or several of these codons as stops. However, this mistake, as far as I understand is not very critical because in my understanding of the authors analyses the most important part is the location of ATGs and where uORFs end affects only their classification as overlapping. Note, that all uORFs in *C. magnum* and other species with no internal stops are oORFs.

Relevant to this is the definition of uORFs provided in the revised version with which I cannot fully agree: "Based on the position and frame of a uORF relative to the downstream CDS, uORFs can be classified into nonoverlapping uORFs, out-of-frame overlapping uORFs (oORFs), and N-terminal extensions"

I strongly disagree with referring to N-terminal extensions as uORFs. If there is an upstream AUG that is in-frame with annotated AUG and initiation takes place on it, we are dealing with misannotated CDS, rather than a distinct phenomenon such as translation of short ORFs occurring upstream of CDS. More than one start could be used for initiation at protein coding ORFs (see our recent work related to this phenomenon in human mRNAs - Benitez-Cantos et al 2020 <https://doi.org/10.1101/gr.257352.119>) and we could refer to the products of these alternative initiation events as truncations or extensions relative to each other, but certainly the longer extended proteoform is not an uORF, it is CDS.

The authors noted that O/E ratio in *C. magnum* is not below 1. Could this be simply because what the authors refer to as uORFs are either CDS starts or internal methionine codons? Hence there is no selection against them? But if it is the case, we get back to the thesis that I stated earlier –

there are no uORFs (in classical sense) in the species with no internal stops.

It seems to me that protists are a special case and their comprehensive analysis could be difficult and inappropriate for this manuscript which is already quite substantial. To some extent the authors already acknowledge this by saying "how the genetic code reassignments affect the distribution and evolution of uORFs in certain protists deserves further study." I think this should be stated more clearly, that the occurrence of uORFs in some protists with genetic code variants may differ substantially from that of most other eukaryotic organisms.

Finally, in relation to the discovery of "stopless" genetic code in *C. magnum* the authors cite one work, however, this discovery was made in two laboratories independently at the same time so both references should be used:

1. Swart et al 2016 <https://doi.org/10.1016/j.cell.2016.06.020>

2. Heaphy et al 2016 <https://doi.org/10.1093/molbev/msw166>

Perhaps the authors may wish to mention other species with the genetic codes which are incompatible with uORFs existence, though this is not necessary unless the authors wish to discuss it:

Blastocrithidia, see Zahonova et al 2016 <https://doi.org/10.1016/j.cub.2016.06.064>)

Ciliate Euplotes, see Lobanov et al 2017 <https://doi.org/10.1038/nsmb.33300>

Amoebophrya sp. ex *Karlodinium veneficum*, see Bachvaroff 2019 <https://doi.org/10.1371/journal.pone.0212912>.

//Pasha Baranov//

Reviewer #2:

Remarks to the Author:

The authors have addressed all my comments. The new data and discussions on canonical and non-noncanonical are very interesting. I would be happy to see this work published.

If I may I want to bring up a few minor points:

1. Figure 1a, middle panel. Almost all species except for human, mouse and fruit fly contains a large fraction of mRNAs without 5' UTR. It is quite unusual, since majority of translation rely on a scanning process in 5' UTR to start translation. Is this due to a lack of accurate annotations on 5' UTR?

2. Supplementary Figure 2. It is interesting that the ATG-like triplets (e.g. ATT, TTG, ATC, GTG) are over-represented in 5' UTR, compared with 3' UTR. Does that indicate that non-canonical uORFs are more preferred within 5' UTR? Or alternatively, is this a consequence of depletion of canonical uORFs, since a single mutation on ATG triplets can easily lead to a switch of canonical uORFs to non-canonical uORFs (ATG -> TTG). As shown in this study, non-canonical uORFs are less effective to repress CDS translation, therefore, these single mutations are not removed during evolution.

3. Rebuttal letter, line 466. "we found that approximately 70-90% of canonical uORFs can be translated (Supplementary Table 3)." It is an interesting data, but I can not find this dataset in Supp. Table 3, or am I missing something? By the way, how to define a translatable uORF, by

using Ribo-seq or TIS-seq. Was that done based on a threshold or by using other ORF prediction tools (e.g. PMID: 26657557)?

Thanks for the opportunity to review this article

Point-to-point response

Reviewer #1 (Remarks to the Author)

The authors performed several analyses which expanded the manuscript substantially and well beyond my expectations. Because of such an extensive revision, a few issues emerged that are specific for the revised version, which I mention below. However, the significance of these issues is relatively minor. I think that the general message of this manuscript is adequately supported.

Response: Thank you for the thorough review of our manuscript. We appreciate the positive feedbacks from this reviewer. We have addressed your comments and concerns in this revised version. Please refer to the point-to-point responses for details.

Comments pertinent for the revised version.

My very first comment related to the apparent impossibility of uORF (if we define uORFs as short translated sequences upstream of CDS) existence in some protists where genetic codes do not allow translation termination of translation in the internal positions of mRNA. In these organisms any translation initiation event would result in a production of a long protein. So I suggested that the authors should simply acknowledge this fact and modify their title and discussion to make it more specific, i.e. plants and animals, but not eukaryotes in general. Instead of following this simple suggestion, the authors chose a hard way, they decided to expand their analysis to include fungi and protists. Of course, this broadens the manuscript and I agree that the use of eukaryotes in the title is appropriate now. However, it made the work not only broader, but also more complicated. While the situation with fungi is similar to that of plants and animals, protists seem to be different and I would like to ask authors to make additional changes in their manuscript to clarify the situation.

Response: Thanks for the comments and suggestions. The major critique raised by this reviewer is that in our previously submitted version, three of the 23 protists do not use the standard genetic code, which makes the identification of uORFs questionable in these three species.

In this revised version, we only considered the 20 protists that use the standard genetic code, and excluded *Condylostoma magnum* and *Parduczia* sp., both of which had no dedicated stop codons (Heaphy et al., 2016; Swart et al., 2016), and *Ichthyophthirius multifiliis*, in which UAA and UAG are reassigned to encode glutamine (Coyne et al., 2011).

We inserted the sentence “Since many protists use alternative nuclear genetic codes involving stop-codon reassignments⁶⁸⁻⁷³ or obligatory frameshifting at internal stop codons⁷⁴, here we only focused on 20 protists that use the standard genetic code (Supplementary Table 1).” in the revised manuscript to clarify this point (Lines 87- 89, Page 4).

First, protists are hugely diverse phylogenetically and exhibit considerable diversity in organisation of their genetic information in their genomes and probably in mRNAs as well. Because of that the analysis of 23 species may not be adequate to obtain a general picture.

Response: Thanks for raising this concern. We emphasized this point in the Discussion with the following sentences (Lines 475-481, Page 15):

“Protists have a very high phylogenetic diversity¹²⁸, and many protists use alternative nuclear genetic codes involving stop-codon reassignments^{68,69} and obligatory frameshifting at internal stop codons⁷⁴. In protists with no dedicated stop codons⁷¹, such as *Condylostoma magnum*^{70,71}, *Parduczia* sp.⁷¹, *Blastocrithidia*⁷², and *Amoebophrya* sp. ex *Karlodinium veneficum*⁷³, translation from any possible uAUG is supposed to terminate

near the end of a transcript and overlaps with the main CDS, which results in a different protein. Thus, the occurrence of uORFs in protists with alternative genetic decoding schemes might differ considerably from that of most other eukaryotes. In this study, we only focused on 20 protists that use the standard genetic code.”

Second, the analysis was done incorrectly. The authors did not take into account the diversity of their genetic codes in these species: “Putative uORFs that start with AUG codons and end with stop codons (UAA/UAG/UGA) were identified from the annotated 5' UTRs of protein-coding genes”. This is clearly wrong for the organisms that do not use one or several of these codons as stops. However, this mistake, as far as I understand is not very critical because in my understanding of the authors analyses the most important part is the location of ATGs and where uORFs end affects only their classification as overlapping. Note, that all uORFs in *C. magnum* and other species with no internal stops are oORFs.

Response: Thanks for pointing this out. In this revised version, we corrected this mistake by only focusing on 20 protists that use the standard genetic code. Please refer to our response above.

Relevant to this is the definition of uORFs provided in the revised version with which I cannot fully agree: “Based on the position and frame of a uORF relative to the downstream CDS, uORFs can be classified into nonoverlapping uORFs, out-of-frame overlapping uORFs (oORFs), and N-terminal extensions”

I strongly disagree with referring to N-terminal extensions as uORFs. If there is an upstream AUG that is in-frame with annotated AUG and initiation takes place on it, we are dealing with misannotated CDS, rather than a distinct phenomenon such as translation of short ORFs occurring upstream of CDS. More than one start could be used for initiation at protein coding ORFs (see our recent work related to this phenomenon in human mRNAs - Benitez-Cantos et al 2020 <https://doi.org/10.1101/gr.257352.119>) and we could refer to the products of these alternative initiation events as truncations or extensions relative to each other, but certainly the longer extended proteoform is not an uORF, it is CDS.

Response: Thank you very much for pointing this out. The nomenclature of ORFs has been inconsistent among different studies, as discussed in a recent review (Orr et al., 2019). For an ORF that has the AUG start codon located in the 5' UTR (defined as “uAUG” in our study), it can function as the start codon of a uORF that has a stop codon either preceding the start codon of the downstream CDS (nonoverlapping uORF, nORF) or residing in the body of the downstream CDS (out-of-frame overlapping uORF, oORF). Also, an uAUG can function as the start codon of an ORF whose stop codon overlaps with the stop codon of the downstream CDS (N-terminal extension, NTE).

Although several studies used a strict definition and only considered nORFs as uORFs (Calviello, et al. 2016; Johnstone, et al. 2016; Whiffin, et al. 2020), many studies broadly treated all the three categories of ORFs as uORFs as they only required the start codons to reside in 5' UTRs (Brar, et al. 2012; Aspden, et al. 2014; Chew, et al. 2016; McGillivray, et al. 2018; Niu, et al. 2020). Several studies also argued that only oORFs and nORFs should be treated as uORFs, and NTEs should be treated as alternative initiation of CDS (Calvo, et al. 2009; Benitez-Cantos, et al. 2020; Chen, et al. 2020).

Previously, we took the broad definition of uORFs and considered all the three categories of uAUGs in 5' UTRs as uORFs. However, as suggested by this reviewer, NTEs are alternative initiation sites of CDSs and thus might differ from nORFs and oORFs in function and sequence evolution. Therefore, in this revision, we followed the suggestions of this reviewer and thoroughly revised the manuscript. In short, we only treated nORFs and oORFs as uORFs in this analysis, although we still considered the start codons of NTEs as a type of uAUGs in this new version. Specifically, we made the following changes:

First, we defined uORFs in the Introduction of the revised manuscript with the following sentences (Line 40-45, Page 3):

“For an AUG triplet in the 5' UTR (defined as “uAUG” hereafter), it can function as the start codon of a uORF that has a stop codon either preceding the start codon of the downstream CDS (nonoverlapping uORF, nORF) or residing in the body of the downstream CDS (out-of-frame overlapping uORF, oORF)^{4,11-18}. Less frequently, an uAUG can function as the start codon of an ORF whose stop codon overlaps with the stop codon of the downstream CDS (N-terminal extension, NTE)^{4,19-21}.”

Second, we also described the compositions of uORFs and NTEs in the examined eukaryotes with the following sentences (Lines 109-112, Page 4):

“The vast majority (> 97%) of the uAUGs identified in the 478 eukaryotic species were start codons of putative canonical uORFs. Specifically, in a species, the percentage (mean \pm s.e.) of nORFs, oORFs, and NTEs was $83.45 \pm 0.41\%$, $14.24 \pm 0.34\%$, and $2.31 \pm 0.15\%$, respectively. The detailed information for the uORFs (nORFs and oORF) and NTEs is presented in Supplementary Table 1.”

Third, we focused our main analyses on the putative canonical uORFs. We emphasized this point with the following sentences (Lines 153-156, Page 6):

“Overall, these results suggest that uAUGs were selected against in 5' UTRs, and the NTEs, which only accounted for a small fraction ($\sim 2.31\%$ on average) of the uAUGs, were also shaped by strong purifying selection during evolution. Since uORFs (nORFs and oORFs) and NTEs might have different mechanisms in regulating gene expression and function, in what follows, we only focused on the putative canonical uORFs.”

Fourth, with the new definition of uORFs, we updated all the relevant analyses that were presented in the main figures (Figs. 1, 2, 3, 5, 6) and supplementary information (Supplementary Figs. 5, 7, 8, 10, 15, and 16; Supplementary Tables 2, 3, 4, 7, and 9). Figures 4 and 7 were not affected since NTEs had been excluded in the previous version.

Despite the extensive updates of the analyses, our conclusions were not affected. This is expected since the NTEs only account for a minor fraction ($2.31 \pm 0.14\%$, mean \pm s.e.) of uAUGs in the species we investigated.

We hope these changes are satisfactory to this reviewer. We are certainly willing to make further revisions if this reviewer thinks additional changes are needed.

The authors noted that O/E ratio in *C. magnum* is not below 1. Could this be simply because what the authors refer to as uORFs are either CDS starts or internal methionine codons? Hence there is no selection against them? But if it is the case, we get back to the thesis that I stated earlier – there are no uORFs (in classical sense) in the species with no internal stops.

Response: Thanks for raising this concern. In this revised version, we removed *C. magnum*, as well as *Parduezia* sp. and *Ichthyophthirius multifiliis*, from the analysis. Please refer to our response above.

However, among the 20 protists that use the standard genetic code, we still found the O/E ratio for uAUGs was close to or higher than 1 in five protists. Nevertheless, we think this observation might be an artifact caused by inaccurate 5' UTR annotations in these five species, because these five protists tended to have significantly longer 5' UTRs than the other 15 protists (Supplementary Fig. 18). Importantly, the O/E ratio of uAUGs in the 5' UTR regions that are proximal to CDS (within 100 nt or 150 nt) were significantly lower than 1 in all the five

protists (Supplementary Table 11), suggesting that uAUG occurrence in 5' UTR regions proximal to CDSs is still under purifying selection in these five protists.

We reported these observations in Discussion with the following sentences (Lines 481-492, Page 15):

“In this study, we only focused on 20 protists that use standard genetic code. Although the O/E ratio of uAUGs was significantly less than 1 in all the fungi, multi-cellular plants and animals we examined, such a pattern was observed in only 15 of the 20 protists. The O/E ratio of uAUGs was close to or higher than 1 in the remaining five protists, including *Cystoisospora suis* (1.161, 95% CI 1.154~1.169), *Toxoplasma gondii* (0.998, 95% CI 0.989~1.1.007), *Nannochloropsis gaditana* (0.997, 95% CI 0.986~1.007), and two malaria vectors *Plasmodium yoelii* (1.016, 95% CI 1.008~1.025) and *Plasmodium vivax* (0.989, 95% CI 0.975~1.004). However, these five protists tended to have significantly longer 5' UTRs than the other 15 protists (Supplementary Fig. 18), suggesting this observation might be an artifact caused by inaccurate 5' UTR annotations in these five species. Indeed, the O/E ratio of uAUGs in the 5' UTR regions that are proximal to CDS (within 100 nt or 150 nt) were significantly lower than 1 in all the five protists (Supplementary Table 11), suggesting that uAUG occurrence in 5' UTR regions proximal to CDSs is still under purifying selection in these protists.”

It seems to me that protists are a special case and their comprehensive analysis could be difficult and inappropriate for this manuscript which is already quite substantial. To some extent the authors already acknowledge this by saying “how the genetic code reassignments affect the distribution and evolution of uORFs in certain protists deserves further study.” I think this should be stated more clearly, that the occurrence of uORFs in some protists with genetic code variants may differ substantially from that of most other eukaryotic organisms.

Response: Thanks for the advice. We rephrased the sentences “Thus, the occurrence of uORFs in protists with alternative genetic decoding schemes might differ considerably from that of most other eukaryotes. In this study, we only focused on 20 protists that use standard genetic code.” to make this point clearer in Discussion (Lines 479-481, Page 15).

Finally, in relation to the discovery of “stopless” genetic code in *C. magnum* the authors cite one work, however, this discovery was made in two laboratories independently at the same time so both references should be used:

1. Swart et al 2016 <https://doi.org/10.1016/j.cell.2016.06.020>
2. Heaphy et al 2016 <https://doi.org/10.1093/molbev/msw166>

Response: Thank you for pointing this out. The citations have been updated (Refs. 70 and 71) in this revision.

Perhaps the authors may wish to mention other species with the genetic codes which are incompatible with uORFs existence, though this is not necessary unless the authors wish to discuss it:

Blastocrithidia, see Zahonova et al 2016 <https://doi.org/10.1016/j.cub.2016.06.064>) Ciliate Euplotes, see Lobanov et al 2017 <https://doi.org/10.1038/nsmb.3330> Amoebophrya sp. ex *Karlodinium veneficum*, see Bachvaroff 2019 <https://doi.org/10.1371/journal.pone.0212912>.

Response: Thank you for these suggestions. These works have been briefly mentioned in Discussion as follows (Lines 475-481, Page 15):

“Protists have a very high phylogenetic diversity¹²⁸, and many protists use alternative nuclear genetic codes involving stop-codon reassignments^{68,69} and obligatory frameshifting at internal stop codons⁷⁴. In protists with no dedicated stop codons⁷¹, such as *Condylostoma magnum*^{70,71}, *Parduczia* sp.⁷¹, *Blastocrithidia*⁷², and *Amoebophrya* sp. ex *Karlodinium veneficum*⁷³, translation from any possible uAUG is supposed to terminate near the end of a transcript and overlaps with the main CDS, which results in a different protein. Thus, the

occurrence of uORFs in protists with alternative genetic decoding schemes might differ considerably from that of most other eukaryotes.”

Reviewer #2 (Remarks to the Author):

The authors have addressed all my comments. The new data and discussions on canonical and non-noncanonical are very interesting. I would be happy to see this work published.

Response: We thank this reviewer for the positive review.

If I may I want to bring up a few minor points:

1. Figure 1a, middle panel. Almost all species except for human, mouse and fruit fly contains a large fraction of mRNAs without 5' UTR. It is quite unusual, since majority of translation rely on a scanning process in 5' UTR to start translation. Is this due to a lack of accurate annotations on 5' UTR?

Response: Thanks for raising this concern. We agree with this reviewer that many non-model organisms lack accurate annotations of 5' UTRs. We rephrased the related sentence (Lines 115-117, Page 5), which reads as follows:

“The number of uAUGs varied wildly across species, either due to the differences in the sequencing coverage of genomes, the accuracy and completeness of 5' UTR annotation, the number of protein-coding genes, the length of 5' UTRs, or mutational bias in 5' UTRs.”

We also emphasized this point in the legend of Fig. 1a with the following sentence (Lines 1,078-1,079, Page 28): “The unavailability of annotated 5' UTRs for many genes in less-studied organisms is presumably caused by the lack of accurate annotations.”

2. Supplementary Figure 2. It is interesting that the ATG-like triplets (e.g. ATT, TTG, ATC, GTG) are over-represented in 5' UTR, compared with 3' UTR. Does that indicate that non-canonical uORFs are more preferred within 5' UTR? Or alternatively, is this a consequence of depletion of canonical uORFs, since a single mutation on ATG triplets can easily lead to a switch of canonical uORFs to non-canonical uORFs (ATG -> TTG). As shown in this study, non-canonical uORFs are less effective to repress CDS translation, therefore, these single mutations are not removed during evolution.

Response: These comments are enlightening. In this revised version, we discussed these two possibilities with the following sentences (Lines 129-134, Page 5):

“Interestingly, some AUG-like triplets (e.g., AUU, UUG, AUC, and GUG) tended to have higher O/E ratios in 5' UTRs than in 3' UTRs in all the clades. Such AUG-like triplets were either selectively maintained in 5' UTRs as they can be used as noncanonical start codons, or alternatively, were the consequence of the depletion of uAUGs because point mutations can easily convert AUG to AUG-like triplets (e.g., from AUG → UUG) in the 5' UTRs. However, further studies are required to separate these two possibilities.”

3. Rebuttal letter, line 466. “we found that approximately 70-90% of canonical uORFs can be translated (Supplementary Table 3).” It is an interesting data, but I can not find this dataset in Supp. Table 3, or am I missing something? By the way, how to define a translatable uORF, by using Ribo-seq or TIS-seq. Was that done based on a threshold or by using other ORF prediction tools (e.g. PMID: 26657557)?

Response: We apologize for the typo. The summary statistics of translated uORFs were presented in Supplementary Table 4, not Supplementary Table 3. We didn't use any ORF prediction tools to define the translatable uORFs. We defined a translated uORF based on a threshold of ribosome-protected fragments (RPFs) whose P-sites are located in this uORF.

In this revised version, we clarified this point in the Methods section with the following sentence (Lines 570-571, Page 17): "A uORF was considered as translated if it was covered by the P-site of at least one RPF read across different ribosome profiling datasets in a species."

We also provided the list of translated uORFs and associated RPF counts in the Source Data of this revised manuscript (figshare doi: [10.6084/m9.figshare.12612068](https://doi.org/10.6084/m9.figshare.12612068)).

References

- Aspden, J.L., Eyre-Walker, Y.C., Phillips, R.J., Amin, U., Mumtaz, M.A.S., Brocard, M., and Couso, J.-P. (2014). Extensive translation of small Open Reading Frames revealed by Poly-Ribo-Seq. *eLife* 3, e03528.
- Bachvaroff, T.R. (2019). A precedented nuclear genetic code with all three termination codons reassigned as sense codons in the syndinean *Amoebophrya* sp. ex *Karlodinium veneficum*. *PLoS One* 14, e0212912.
- Baranov, P.V., Atkins, J.F., and Yordanova, M.M. (2015). Augmented genetic decoding: global, local and temporal alterations of decoding processes and codon meaning. *Nat Rev Genet* 16, 517-529.
- Brar, G.A., Yassour, M., Friedman, N., Regev, A., Ingolia, N.T., and Weissman, J.S. (2012). High-resolution view of the yeast meiotic program revealed by ribosome profiling. *Science* 335, 552-557.
- Burki, F., Roger, A.J., Brown, M.W., and Simpson, A.G.B. (2020). The New Tree of Eukaryotes. *Trends Ecol Evol* 35, 43-55.
- Calviello, L., Mukherjee, N., Wyler, E., Zauber, H., Hirsekorn, A., Selbach, M., Landthaler, M., Obermayer, B., and Ohler, U. (2016). Detecting actively translated open reading frames in ribosome profiling data. *Nat Methods* 13, 165-170.
- Calvo, S.E., Pagliarini, D.J., and Mootha, V.K. (2009). Upstream open reading frames cause widespread reduction of protein expression and are polymorphic among humans. *Proc Natl Acad Sci U S A* 106, 7507-7512.
- Chen, J., Brunner, A.D., Cogan, J.Z., Nuñez, J.K., Fields, A.P., Adamson, B., Itzhak, D.N., Li, J.Y., Mann, M., Leonetti, M.D., *et al.* (2020). Pervasive functional translation of noncanonical human open reading frames. *Science* 367, 1140-1146.
- Chew, G.L., Pauli, A., and Schier, A.F. (2016). Conservation of uORF repressiveness and sequence features in mouse, human and zebrafish. *Nat Commun* 7, 11663.
- Coyne, R.S., Hannick, L., Shanmugam, D., Hostetler, J.B., Brami, D., Joardar, V.S., Johnson, J., Radune, D., Singh, I., Badger, J.H., *et al.* (2011). Comparative genomics of the pathogenic ciliate *Ichthyophthirius multifiliis*, its free-living relatives and a host species provide insights into adoption of a parasitic lifestyle and prospects for disease control. *Genome Biol* 12, R100.
- Heaphy, S.M., Mariotti, M., Gladyshev, V.N., Atkins, J.F., and Baranov, P.V. (2016). Novel Ciliate Genetic Code Variants Including the Reassignment of All Three Stop Codons to Sense Codons in *Condylostoma magnum*. *Molecular biology and evolution* 33, 2885-2889.
- Johnstone, T.G., Bazzini, A.A., and Giraldez, A.J. (2016). Upstream ORFs are prevalent translational repressors in vertebrates. *The EMBO journal* 35, 706-723.
- Lobanov, A.V., Heaphy, S.M., Turanov, A.A., Gerashchenko, M.V., Pucciarelli, S., Devaraj, R.R., Xie, F., Petyuk, V.A., Smith, R.D., Klobutcher, L.A., *et al.* (2017). Position-dependent termination and widespread obligatory frameshifting in *Euplotes* translation. *Nat Struct Mol Biol* 24, 61-68.

- McGillivray, P., Ault, R., Pawashe, M., Kitchen, R., Balasubramanian, S., and Gerstein, M. (2018). A comprehensive catalog of predicted functional upstream open reading frames in humans. *Nucleic acids research* *46*, 3326-3338.
- Niu, R., Zhou, Y., Zhang, Y., Mou, R., Tang, Z., Wang, Z., Zhou, G., Guo, S., Yuan, M., and Xu, G. (2020). uORFlight: a vehicle toward uORF-mediated translational regulation mechanisms in eukaryotes. *Database (Oxford)* *2020*.
- Orr, M.W., Mao, Y., Storz, G., and Qian, S.-B. (2019). Alternative ORFs and small ORFs: shedding light on the dark proteome. *Nucleic Acids Research*.
- Sengupta, S., and Higgs, P.G. (2015). Pathways of Genetic Code Evolution in Ancient and Modern Organisms. *J Mol Evol* *80*, 229-243.
- Swart, E.C., Serra, V., Petroni, G., and Nowacki, M. (2016). Genetic Codes with No Dedicated Stop Codon: Context-Dependent Translation Termination. *Cell* *166*, 691-702.
- Whiffin, N., Karczewski, K.J., Zhang, X., Chothani, S., Smith, M.J., Evans, D.G., Roberts, A.M., Quaipe, N.M., Schafer, S., Rackham, O., *et al.* (2020). Characterising the loss-of-function impact of 5' untranslated region variants in 15,708 individuals. *Nat Commun* *11*, 2523.
- Záhonová, K., Kostygov, A.Y., Ševčíková, T., Yurchenko, V., and Eliáš, M. (2016). An Unprecedented Non-canonical Nuclear Genetic Code with All Three Termination Codons Reassigned as Sense Codons. *Curr Biol* *26*, 2364-2369.

Reviewers' Comments:

Reviewer #1:

Remarks to the Author:

The authors addressed all my comments in full. I would like to congratulate the authors on the comprehensive and timely study dedicated to the important topic of the uORFs evolution.

Pavel Baranov.

Point-to-point response

Reviewer #1 (Remarks to the Author):

The authors addressed all my comments in full. I would like to congratulate the authors on the comprehensive and timely study dedicated to the important topic of the uORFs evolution.

Pavel Baranov.

[Response: We thank this reviewer for the positive feedback.](#)