

1
2 A spectral analysis approach to detect actively translated open reading
3 frames in high-resolution ribosome profiling data
4
5

6 Lorenzo Calviello¹, Neelanjan Mukherjee^{1*}, Emanuel Wyler^{1*}, Henrik Zauber¹, Antje
7 Hirse Korn¹, Matthias Selbach¹, Markus Landthaler¹, Benedikt Obermayer¹, Uwe Ohler^{1,2}

8
9 ¹Berlin Institute for Medical Systems Biology, Max Delbrueck Center for Molecular
10 Medicine, 13125 Berlin

11 ²Departments of Biology and Computer Science, Humboldt University, 10099 Berlin
12

13
14 *equal contribution. Corresponding author: uwe.ohler@mdc-berlin.de

15
16
17 Abstract
18

19 RNA sequencing protocols allow for quantifying gene expression regulation at each individual step, from
20 transcription to protein synthesis. Ribosome Profiling (Ribo-seq) maps the positions of translating
21 ribosomes over the entire transcriptome. Despite its great potential, a rigorous statistical approach to
22 identify translated regions by means of the characteristic three-nucleotide periodicity of Ribo-seq data is
23 not yet available. To fill this gap, we developed RiboTaper, which quantifies the significance of periodic
24 Ribo-seq reads via spectral analysis methods.

25 We applied RiboTaper on newly generated, deep Ribo-seq data in HEK293 cells, to derive an extensive
26 map of translation that covers Open Reading Frame (ORF) annotations for more than 11,000 protein-
27 coding genes. We also find distinct ribosomal signatures for several hundred detected upstream ORFs
28 and ORFs in annotated non-coding genes (ncORFs). Mass spectrometry data confirms that RiboTaper
29 achieves excellent coverage of the cellular proteome and validates dozens of novel peptide products.
30 Collectively, RiboTaper (available at <https://ohlerlab.mdc-berlin.de/software/>) is a powerful method for
31 comprehensive *de novo* identification of actively used ORFs in the human genome.
32
33
34
35

1 INTRODUCTION

2

3 Ribosome Profiling (Ribo-seq) is a method to globally investigate protein synthesis¹ by sequencing RNA
4 fragments protected by engaged ribosomes. Mapping these ribosomal footprints to the transcriptome
5 produces a detailed quantitative picture of translation across thousands of genes². In the few years since
6 its first description, Ribo-seq has been applied to monitor translation in different organisms³ and
7 biological conditions⁴. Applications cover a wide range of subjects, including the identification of
8 alternative start codons⁵ and the definition of short peptides^{6,7} as well as quantitative aspects such as
9 translation rates or codon efficiency⁸.

10

11 With an increasing number of annotated non-coding RNAs, a central question is whether some of these
12 transcripts contain short translated ORFs missed by annotation efforts. To this end, a number of studies
13 proposed *ad-hoc* Ribo-seq based global metrics that aimed at defining the translational status of such
14 non-coding transcripts⁹⁻¹¹. The mere presence of Ribo-seq reads in regions of the transcriptome does
15 however not imply the presence of actively elongating ribosomes. The mode of Ribo-seq read length
16 distributions typically falls at ~29-30nt, which is the known fragment size protected by 80S ribosomes¹².
17 Consequently, this subset of ribosomal footprints displays a striking bias towards the translated frame
18 ^{1,6,9}, which can be used to infer, for each read, the position of the peptidyl-site (P-site) compartment of
19 the translating ribosome⁹. This sub-codon resolution holds great promise to discriminate between
20 ribosomal coverage and a periodic footprint profile (PFP), i.e., a consistent, 3-periodic codon-by-codon
21 alignment pattern across a transcript. Yet, this property has only recently been exploited within
22 summary statistics to identify small translated ORFs^{6,13} or to explore translation on multiple frames¹⁴.

23

24 Despite the recent development Ribo-seq analysis tools¹⁵⁻¹⁷, a statistically rigorous method using this
25 property to identify translated regions has not been proposed. Here we present a computational
26 approach, based on spectral analysis, to comprehensively identify the set of PFPs in a given Ribo-seq
27 sample. Our method, called RiboTaper, exploits the sub-codon resolution of Ribo-seq reads to call high-
28 confidence translated loci, and reconstructs the full set of ORFs in annotated coding and non-coding
29 transcripts. We quantify how RiboTaper outperforms alternative approaches, show that a limited
30 number of non-coding RNAs contain actively translated ORFs, analyze evolutionary signatures in
31 different ORF classes, and cross-reference the identified ORFs with proteome-wide experimental
32 evidence.

1

2 RESULTS

3

4 *The RiboTaper strategy for testing Ribo-seq sub-codon profiles*

5

6 The main idea of our approach is to provide a statistical test to identify PFPs, indicative of consistent
7 codon-by-codon ribosomal movement across a putative translated region. This test is the central part of
8 the RiboTaper method (Fig. 1a): We first define P-site positions for the majority of mapped reads,
9 according to their aggregate profiles over annotated start codons^{6,9} (Fig. 1a, Supplementary Fig. 1). Next,
10 we create data tracks for every annotated exonic region, and use the multitaper approach¹⁸ to test for
11 significant PFPs in exonic P-site profiles (Fig. 1a). Finally, exonic profiles are merged according to the
12 annotated transcript structures to detect translation on *de novo* identified ORFs (cf. Methods, Software
13 availability, Supplementary Software).

14

15 The multitaper approach applies a set of orthogonal window functions (tapers) to a discrete signal
16 before computing its Fourier transform. The transformed windowed outputs are averaged to obtain a
17 smoothed spectrum amenable to a non-parametric test for detecting significant frequencies^{19,20} (cf.
18 Methods). In this way, the multitaper method tests directly for the presence of PFPs, in contrast to
19 current approaches that look at enrichment of Ribo-seq reads over RNA-seq (Translation Efficiency¹) or a
20 preference for reads to align to one frame.

21

22 *RiboTaper defines active translation with high sensitivity and specificity*

23

24 To assess the performance and utility of RiboTaper, we generated a deep ribosome profiling dataset in
25 HEK293 cells following established protocols (cf. Methods). We obtained >29 Mio reads, out of which
26 >25 Mio aligned to the genome (Supplementary Table 1). Profiles at consensus coding exons (CCDS
27 annotation, see Methods) displayed a striking frame preference, in excellent agreement with annotated
28 transcript types and regions (Supplementary Fig. 2).

29

30 To evaluate the sensitivity of the multitaper, we applied it on profiles from CCDS exons of different
31 length and Ribo-seq coverage (Supplementary Fig. 2-4). As expected, the test achieved higher sensitivity

1 on longer exons and benefitted from higher coverage (Fig 1b). Twenty-four tapers exhibited the best
2 combination of sensitivity and specificity (Supplementary Fig. 5).

3
4 We then benchmarked the multitaper against a Chi-squared significance test as baseline. This test uses a
5 null hypotheses of a uniform frame P-site distribution and corresponds to the basic assumption behind
6 the ORFscore⁶. However, sequencing protocols are affected by different sources of variability that cause
7 non-uniform distribution of reads²¹. Further, insufficient depth may also lead to sampling biases and
8 spurious enrichments. To evaluate specificity on an appropriate negative sample, i.e. regions without
9 3nt-periodic reads, we applied the tests on RNA-seq data of annotated CCDS exons.

10
11 At a significance level of 0.05, the multitaper displayed slightly lower sensitivity when compared to the
12 Chi-squared test (87% vs. 94% of positive exons, Supplementary Fig. 2-4). However, when applied to
13 RNA-seq profiles, the Chi-squared test showed a worrisome number of positive calls compared to the
14 multitaper test (16% vs. 3.5%, Fig. 1c, Supplementary Fig 2-4). We observed strongly skewed p-values for
15 the Chi-squared test but a desirable near-uniform distribution of the multitaper p-values on the RNA-seq
16 (Fig 1c). This high specificity is critical when exploring translation outside of protein-coding regions and
17 directly pertains to Ribo-seq data: Integrated in our analysis pipeline, the Chi-square detected 67
18 translated ORFs in snoRNAs and snRNAs – transcripts unlikely to generate peptides -- while the
19 multitaper detected only 3. A principled statistical framework eliminates the need for heuristic
20 parameters or ad-hoc cutoffs, a crucial step when applying score-based metrics on different datasets
21 (Supplementary Fig. 6).

22

23 *Full ORF reconstruction captures known and novel translated ORFs*

24

25 We next created transcript tracks for *de novo* translated ORF identification, using 3nt periodicity and
26 frame definition by the P-site positions (Fig 2a, Methods). We classified ORFs based on annotation
27 categories and genomic position relative to known coding regions (Fig. 2b-c, Supplementary Fig. 7,
28 Methods). This led to ~21,000 translated ORFs in ~14,000 expressed genes, in coding and non-coding
29 transcripts, across a wide range of expression values (Fig. 3a).

30

31 The vast majority of ORFs in protein-coding genes overlapped known CCDS coding regions; 369 non-
32 CCDS protein-coding genes were identified as harboring translated ORFs (“nonccds_coding_ORFs”). We

1 detected >600 genes with translated upstream ORFs (uORFs) and 54 genes with downstream ORFs
2 (dORFs; cf. Methods). We additionally identified ORFs in 504 non-coding genes (ncORFs). These ORFs
3 belong mainly to pseudogene, antisense, and lincRNA biotypes (Fig. 3b).

4
5 RiboTaper identified ORFs exhibited a protein-coding like distribution of Ribo-seq read lengths as
6 quantified by the FLOSS score¹¹, across all ORF categories (Supplementary Fig. 8). The reconstructed ORF
7 coordinates agreed with translation initiation sites defined by QTI-seq²² or the reference annotation (Fig.
8 3c). Compared with the reference, 149 upstream initiation sites were detected by both QTI-seq and
9 RiboTaper, mostly corresponding to uORF start codons (Fig. 3c). 52 internal starts were identified by
10 both QTI-seq and our method. Approximately 1000 QTI-seq ATG start codon candidates did not overlap
11 with either annotated or RiboTaper-defined start codons. Using Ribo-seq data from the same study,
12 more than 99% of RiboTaper-identified CCDS ORFs were also found in our data (Fig. 3d). Agreement
13 dropped to 68% for lincRNAs/antisense ORFs and 47% for uORF-containing genes, possibly due to their
14 relatively short length and low expression levels (Fig. 3d).

15
16 To demonstrate the general applicability of RiboTaper, we identified thousands of coding ORFs and
17 ncORFs in Ribo-seq data of the zebrafish embryo⁶ (Supplementary Fig. 4, Supplementary Table 2-3,
18 Supplementary File 1). Among the identified ncORFs was the recently discovered ORF in the lincRNA
19 *toddler*²³, which encodes a small polypeptide morphogen essential for zebrafish embryonic development
20 (Supplementary Fig. 9).

21

22

23 *Conservation patterns underline diverse roles of open reading frames in different* 24 *RNA classes*

25

26 To gain insights into the identified ORFs in categories outside of annotated coding regions, we analyzed
27 evolutionary conservation patterns of ORFs in different categories. CCDS ORFs and non-CCDS coding
28 ORFs exhibited high nucleotide conservation²⁴, followed by pseudogenes and processed transcripts
29 (Supplementary Fig. 10). Different from *bona fide* coding regions, we observed elevated nucleotide
30 conservation only around start and stop codons for uORFs and processed transcripts (Fig. 4a).

31

1 Next, we assessed whether sequence conservation reflected evolutionary selection on the encoded
2 protein sequence, and whether this selection persists in the human population. We quantified coding
3 potential by means of hexamer sequence statistics²⁵ and tested the preference for synonymous vs. non-
4 synonymous SNPs in the human population using appropriate length- and conservation-matched
5 controls(Fig. 4c, Supplementary Fig. 10, Methods). For CCDS and non-CCDS coding, ORFs, as well as
6 processed_transcript ncORFs, nucleotide conservation was accompanied by good hexamer scores and a
7 depletion of nonsynonymous SNPs (dN/dS). uORFs also showed conservation on the nucleotide level but
8 no significant enrichment of synonymous substitutions (dN/dS), suggesting potential regulatory rather
9 than protein-coding roles, since their position but not sequence is conserved.
10 Additional ncORFs categories showed low nucleotide conservation (except for pseudogenes), a positive
11 trend for hexamer scores, but no depletion for nonsynonymous SNPs. Codon substitution patterns
12 across vertebrate species²⁶ led to a similar outcome (Supplementary Fig. 10). Taken together, ncORFs
13 detected by RiboTaper do not necessarily entail conservation and selection patterns similar to protein-
14 coding regions.

15

16 *Ribosome profiling serves as an effective proxy to define the cellular proteome*

17

18 The ORFs identified by RiboTaper covered a wide range of expression values (cf Fig. 3a), to an extent
19 that its coverage might exceed deep mass spectrometry datasets in defining the cellular proteome. To
20 evaluate this, we created a custom database from the set of identified ORFs to match the spectra of a
21 recent HEK293 tandem mass spectrometry dataset²⁷. The RiboTaper peptide set corresponded to ~59%
22 of the peptides in Uniprot (human entries, rel. October 2014), and an additional 2% of non-Uniprot
23 candidates. The RiboTaper set matched >90.000 peptide sequences, belonging to >8.000 genes (Fig. 4d,
24 Supplementary Fig. 11), similar to the results of the full Uniprot search.

25

26 Over 3900 peptide sequences were found only by our custom search but not using the Uniprot database
27 (1% FDR, Supplementary Table 4). In turn, our search missed a similar number of peptides. RiboTaper-
28 only peptides matched more spectra than UniProt-only peptides, despite being shorter and with lower
29 matching scores (Supplementary Fig. 11). We found little evidence for expression or translation for most
30 of the Uniprot-only peptides (Fig. 4e), suggesting that those may derive from erroneous calls or stable
31 peptides from unstable RNAs.

32

1 RiboTaper ORFs with peptide support mapped to CCDS genes, with few exceptions (Fig. 4f). We
2 identified peptides belonging to a uORF in the MIEF1 gene²⁸ (Fig. 2b) or from dORFs and ncORFs located
3 in conserved and non-conserved genomic regions (Fig. 2c, Supplementary Fig. 12-15). In total, 228
4 identified peptide sequences were not annotated in Uniprot. In many cases, the novel identified
5 peptides mapped uniquely to their respective ORFs (Fig. 4f, Supplementary Table 4).

6

7 DISCUSSION

8

9 Ribo-seq reads of specific lengths can display a precise sub-codon pattern, which allows for accurate
10 identification of the translated frame. However, different experimental protocols can have a marked
11 influence on the sub-codon profiles² (Supplementary Fig. 1). The Ribo-seq protocol is far from being
12 standardized²⁹, and codon biases and ribosome stalling³⁰ pose challenges to the quantitative
13 understanding of translation at each locus. The RiboTaper method proposed here is robust and tailored
14 to exploit the periodic sub-codon pattern to identify periodic footprint profiles (PFPs) associated with
15 active translation. This principled statistical approach based on the significance of spectra allowed us to
16 independently test transcript regions for PFPs on single loci, yielding high specificity over a wide range of
17 expression and coverage.

18

19 Despite the identification of many PFPs in non-coding transcripts^{10,11}, evolutionary conservation analysis
20 suggests that the act of translation rather than the production of specific peptides may be the relevant
21 process and may explain the lack of matching peptides³¹. The notable exception was the apparent
22 translation of a considerable number pseudogenes, which may have retained coding potential but are
23 no longer under purifying selection. Quantifying the presence and significance of ribosome footprint
24 reads becomes increasingly difficult for very small translated regions (<20 amino acid long “dwORFs”⁷),
25 and these may be missed by both RiboTaper as well as by conventional mass-spec protocols.

26

27 Further investigation is needed to understand the function of ribosomal readout at ncORFs and
28 u/dORFs, considering features of the dynamics of translation (frameshifting, reinitiation), together with
29 other aspects of RNA metabolism which are regulated by ribosomal activity. For example, uORFs can
30 promote transcript susceptibility to RNA-surveillance mechanisms³², such as Nonsense-Mediated-Decay.
31 Defining the ensemble of translated uORFs from high-throughput experiments has remained elusive³³,
32 and our method holds promise in the identification of such events.

1

2 Different studies reported a widespread presence of pile-ups at canonical and non-canonical start-
3 codons (NUG) in 5'UTRs²². We here considered only AUG start codons, likely missing cases of NUG
4 starting ORFs (Supplementary Fig. 13). Further analyses are needed to investigate the validity of NUGs as
5 efficient start codons. Another possible extension of our current method is the definition of translated
6 alternative RNA isoforms^{14,34}. Our approach also does not yet account for frameshifting and may miss
7 cases in which multiple actively translated ORFs overlap with each other¹⁴.

8

9 Recent approaches have used Ribo-seq results to aid peptide discovery³⁵. In our hands, the set of
10 identified translated ORFs represented a comparable alternative to public databases and, in fact, a more
11 comprehensive proxy to define the cellular proteome, as the set of RiboTaper PFPs exceeded the
12 coverage of deep mass-spec datasets. Altogether, our approach constitutes a resourceful toolkit for
13 dedicated computational analysis of high-resolution data from Ribo-seq experiments.

14

15

16 Author Contributions:

17 LC and UO developed the computational approach. LC implemented the RiboTaper method and
18 analyzed the sequencing data, supervised by UO. EW, NM and AH performed the Ribo-seq experiments,
19 supervised by ML. BO carried the evolutionary conservation analysis and helped with the interpretation
20 of the presented findings. HZ, MS and LC analyzed the mass spec data. LC, NM and UO wrote the
21 manuscript, with crucial input from all authors.

22

23 Acknowledgments:

24 LC wants to sincerely thank Roberto Marangoni (University of Pisa) for inspiring and fruitful discussions
25 and Alina Munteanu for help with the Supplementary Figures.

26

27 Funding:

28 LC is funded by the MDC PhD program. BO acknowledges funding through a Delbrück fellowship at the
29 MDC. NM acknowledges funding from EU Marie Curie IIF (EU). NM and UO were supported by NIH grant
30 R01GM104962.

31

32 Conflicts of interest:

1 None declared.

2

REFERENCES

- 3 1. Ingolia, N. T., Ghaemmaghami, S., Newman, J. R. S. & Weissman, J. S. Genome-wide analysis in
4 vivo of translation with nucleotide resolution using ribosome profiling. *Science* **324**, 218–223
5 (2009).
- 6 2. Ingolia, N. T., Brar, G. a, Rouskin, S., McGeachy, A. M. & Weissman, J. S. The ribosome profiling
7 strategy for monitoring translation in vivo by deep sequencing of ribosome-protected mRNA
8 fragments. *Nat. Protoc.* **7**, 1534–1550 (2012).
- 9 3. Schafer, S. *et al.* Translational regulation shapes the molecular landscape of complex disease
10 phenotypes. *Nat. Commun.* **6**, 7200 (2015).
- 11 4. Lareau, L. F., Hite, D. H., Hogan, G. J. & Brown, P. O. Distinct stages of the translation elongation
12 cycle revealed by sequencing ribosome-protected mRNA fragments. *Elife* **2014**, 1–16 (2014).
- 13 5. Fritsch, C. *et al.* Genome-wide search for novel human uORFs and N-terminal protein extensions
14 using ribosomal footprinting. *Genome Res.* **22**, 2208–2218 (2012).
- 15 6. Bazzini, A. a. *et al.* Identification of small ORFs in vertebrates using ribosome footprinting and
16 evolutionary conservation. *EMBO J.* **33**, 981–993 (2014).
- 17 7. Aspden, J. L. *et al.* Extensive translation of small Open Reading Frames revealed by Poly-Ribo-Seq.
18 *Elife* **3**, e03528 (2014).
- 19 8. Pop, C. *et al.* Causal signals between codon bias, mRNA structure, and the efficiency of
20 translation and elongation. *Mol. Syst. Biol.* **10**, (2014).
- 21 9. Chew, G.-L. *et al.* Ribosome profiling reveals resemblance between long non-coding RNAs and 5'
22 leaders of coding RNAs. *Development* **140**, 2828–34 (2013).
- 23 10. Guttman, M., Russell, P., Ingolia, N. T., Weissman, J. S. & Lander, E. S. Ribosome profiling provides
24 evidence that large noncoding RNAs do not encode proteins. *Cell* **154**, 240–251 (2013).
- 25 11. Ingolia, N. T. *et al.* Ribosome Profiling Reveals Pervasive Translation Outside of Annotated
26 Protein-Coding Genes. *Cell Rep.* 1365–1379 (2014).
- 27 12. Steitz, J. A. Polypeptide chain initiation: nucleotide sequences of the three ribosomal binding
28 sites in bacteriophage R17 RNA. *Nature* **224**, 957–964 (1969).
- 29 13. Duncan, C. D. S. & Mata, J. The translational landscape of fission-yeast meiosis and sporulation.
30 *Nat. Struct. Mol. Biol.* **21**, 641–7 (2014).

- 1 14. Michel, A. M. *et al.* Observation of dually decoded regions of the human genome using ribosome
2 profiling data. *Genome Res.* **22**, 2219–2229 (2012).
- 3 15. Michel, A. M. *et al.* GWIPS-viz: development of a ribo-seq genome browser. *Nucleic Acids Res.* **42**,
4 D859–64 (2014).
- 5 16. Olshen, A. B. *et al.* Assessing gene-level translational control from ribosome profiling.
6 *Bioinformatics* **29**, 2995–3002 (2013).
- 7 17. Legendre, R., Baudin-Baillieu, A., Hatin, I. & Namy, O. RiboTools: a Galaxy toolbox for qualitative
8 ribosome profiling analysis. *Bioinformatics* **31**, 2586–8 (2015).
- 9 18. Thomson, D. J. Spectrum estimation and harmonic analysis. *Proc. IEEE* **70**, 1055–1096 (1982).
- 10 19. Babadi, B. & Brown, E. N. A review of multitaper spectral analysis. *IEEE Trans. Biomed. Eng.* **61**,
11 1555–64 (2014).
- 12 20. Thomson, D. J., MacLennan, C. G. & Lanzerotti, L. J. Propagation of solar oscillations through the
13 interplanetary medium. *Nature* **376**, 139–144 (1995).
- 14 21. Lahens, N. F. *et al.* IVT-seq reveals extreme bias in RNA-sequencing. *Genome Biol.* **15**, R86 (2014).
- 15 22. Gao, X. *et al.* Quantitative profiling of initiating ribosomes in vivo. *Nat. Methods* **12**, 147–53
16 (2015).
- 17 23. Pauli, A. *et al.* Toddler: an embryonic signal that promotes cell movement via Apelin receptors.
18 *Science* **343**, 1248636 (2014).
- 19 24. Siepel, A. *et al.* Evolutionarily conserved elements in vertebrate, insect, worm, and yeast
20 genomes. *Genome Res.* **15**, 1034–50 (2005).
- 21 25. Wang, L. *et al.* CPAT: Coding-Potential Assessment Tool using an alignment-free logistic
22 regression model. *Nucleic Acids Res.* **41**, e74 (2013).
- 23 26. Lin, M. F., Jungreis, I. & Kellis, M. PhyloCSF: A comparative genomics method to distinguish
24 protein coding and non-coding regions. *Bioinformatics* **27**, 275–282 (2011).
- 25 27. Eravci, M., Sommer, C. & Selbach, M. IPG strip-based peptide fractionation for shotgun
26 proteomics. *Methods Mol. Biol.* **1156**, 67–77 (2014).
- 27 28. Andreev, D. E. *et al.* Translation of 5' leaders is pervasive in genes resistant to eIF2 repression.
28 *Elife* **4**, e03971 (2015).
- 29 29. Gerashchenko, M. V & Gladyshev, V. N. Translation inhibitors cause abnormalities in ribosome
30 profiling experiments. *Nucleic Acids Res.* **42**, 1–7 (2014).

- 1 30. Artieri, C. G. & Fraser, H. B. Accounting for biases in riboprofiling data indicates a major role for
2 proline in stalling translation. *Genome Res.* **24**, 2011–21 (2014).
- 3 31. Jia, H. *et al.* Long noncoding RNAs are rarely translated in two human cell lines. *Genome Res.*
4 1646–1657 (2012).
- 5 32. Barbosa, C., Peixeiro, I. & Romão, L. Gene Expression Regulation by Upstream Open Reading
6 Frames and Human Disease. *PLoS Genet.* **9**, 1–12 (2013).
- 7 33. Wethmar, K., Barbosa-Silva, A., Andrade-Navarro, M. A. & Leutz, A. uORFdb—a comprehensive
8 literature database on eukaryotic uORF biology. *Nucleic Acids Res.* **42**, D60–7 (2014).
- 9 34. Zupanic, A. *et al.* Detecting translational regulation by change point analysis of ribosome profiling
10 data sets. *RNA* **2014**, (2014).
- 11 35. Crappé, J. *et al.* PROTEOFORMER: deep proteome coverage through ribosome profiling and MS
12 integration. *Nucleic Acids Res.* **43**, (2014).
- 13 36. Schueler, M. *et al.* Differential protein occupancy profiling of the mRNA transcriptome. *Genome*
14 *Biol.* **15**, R15 (2014).
- 15 37. Langmead, B., Trapnell, C., Pop, M. & Salzberg, S. L. Ultrafast and memory-efficient alignment of
16 short DNA sequences to the human genome. *Genome Biol.* **10**, R25 (2009).
- 17 38. Dobin, A. *et al.* STAR: ultrafast universal RNA-seq aligner. *Bioinformatics* **29**, 15–21 (2013).
- 18 39. Harrow, J. *et al.* GENCODE: the reference human genome annotation for The ENCODE Project.
19 *Genome Res.* **22**, 1760–74 (2012).
- 20 40. Li, B. & Dewey, C. N. RSEM: accurate transcript quantification from RNA-Seq data with or without
21 a reference genome. *BMC Bioinformatics* **12**, 323 (2011).
- 22 41. Quinlan, A. R. & Hall, I. M. BEDTools: a flexible suite of utilities for comparing genomic features.
23 *Bioinformatics* **26**, 841–2 (2010).
- 24 42. Rahim, K. J., Burr, W. S. & Thomson, D. J. Appendix A: Multitaper R Package in ‘Applications of
25 Multitaper Spectral Analysis to Nonstationary Data,’ PhD diss., Queen’s University,. 149–183
26 (2014). at <<http://hdl.handle.net/1974/12584>>
- 27 43. Mackowiak, S. D. *et al.* Extensive identification and analysis of conserved small ORFs in animals.
28 *Genome Biol.* **16**, 179 (2015).
- 29 44. Cox, J. & Mann, M. MaxQuant enables high peptide identification rates, individualized p.p.b.-
30 range mass accuracies and proteome-wide protein quantification. *Nat. Biotechnol.* **26**, 1367–
31 1372 (2008).

1 45. Chen, C., Li, Z., Huang, H., Suzek, B. E. & Wu, C. H. A fast Peptide Match service for UniProt
2 Knowledgebase. *Bioinformatics* **29**, 2808–9 (2013).

3

4

5 METHODS

6

7 **Ribosome profiling**

8 We followed the original protocol² with minor modifications. For cell lysis, the cell medium was
9 aspirated and cells washed with ice-cold PBS containing 100 µg/ml cycloheximide. No cycloheximide was
10 added to the culture medium before. After thorough removal of the PBS, the plates were immersed in
11 liquid nitrogen and placed on dried ice. For cell lysis, 400 µl mammalian polysome buffer (20 mM Tris-
12 HCl pH 7.4, 150 mM NaCl, 5 mM MgCl₂, with 1 mM DTT and 100 µg/ml cycloheximide added freshly)
13 was supplied with 1% (v/v) Triton X-100 and 25 U/ml Turbo DNase (Life Technologies, AM2238) and
14 dripped on the plate which is subsequently placed tilted on wet ice. The cells were scraped off to the
15 lower portion of the dish so that they thawed in lysis buffer. After dispersal of the cells by pipetting, the
16 lysate was triturated ten times through a 26G needle, cleared by centrifugation at 20000g for 5 minutes,
17 flash-frozen in liquid nitrogen and stored at -80 °C until further usage. For isolation of ribosome-
18 protected RNA fragments, 120 µl of the lysate were digested with 3 µl RNase I (Life Technologies,
19 AM2294) for 45 min. at room temperature with rotation. The digestion was stopped by addition of 4 µl
20 Superase-In (Life Technologies, AM2694). Meanwhile, MicroSpin S-400 HR columns (GE Healthcare, 27-
21 5140-01) were equilibrated with 3 ml mammalian polysome buffer by gravity flow and emptied by
22 centrifugation at 600g for four minutes. Of the digested lysate, 100 µl were then immediately loaded on
23 the column and eluted by centrifugation at 600g for two minutes. RNA was extracted from the flow-
24 through, approximately 125 µl, using Trizol LS (Life Technologies, 10296-010). Ribosomal RNA fragments
25 were then removed using the RiboZero Kit (Illumina, MRZH11124) and separated on a 17% denaturing
26 Urea-PAGE gel (National Diagnostics, EC-829). The size range from 27nt to 30nt as defined by loading 20
27 pmol each Marker-27nt and Marker-30nt was cut out and the RNA fragments subjected to library
28 generation using 3'-Adapter NN-RA3, 5' adapter OR5-NN, RT primer RTP and PCR primers RP1 (forward
29 primer) and RPI6-7 (reverse primer, containing barcodes). Libraries were sequenced on a HiSeq 2000
30 device (Illumina). After initial quality control, we obtained ~29 Mio raw reads by pooling the RPI6 and
31 RPI7 samples.

1
2 Marker-27nt: 5'-AUGUACACGGAGUCGAGCUCAACCCGC-P
3 Marker-30nt: 5'-AUGUACACGGAGUCGAGCUCAACCCGCAAC-P
4 NN-RA3: P NNTGGAATTCTCGGGTGCCAAGG-InvdT
5 OR5-NN: 5'-GUUCAGAGUUCUACAGUCCGACGAUCNN
6 RTP 5'-GCCTTGGCACCCGAGAATTCCA
7 RP1 5'-AATGATACGGCGACCACCGAGATCTACACGTTTCAGAGTTCTACAGTCCGA
8 RPI6 5'-CAAGCAGAAGACGGCATAACGAGATATTGGCGTGACTGGAGTTCCTTGGCACCCGAGAATTCCA
9 RPI7 5'-CAAGCAGAAGACGGCATAACGAGATGATCTGGTGACTGGAGTTCCTTGGCACCCGAGAATTCCA
10

11 **Preprocessing of Ribo-seq and RNA-seq.**

12 Poly-A selected RNA-seq data for HEK293 was obtained from a recent study³⁶ (accession: GSM1306496).
13 Ribo-seq reads were stripped from the adapter sequences, and reads aligning to rRNA sequences were
14 discarded using Bowtie³⁷. Unaligned Ribo-seq reads and RNA-seq reads were aligned to the human
15 genome (hg 19) using the split-aware aligner STAR³⁸. The STAR genome index was built using annotation
16 obtained from GENCODE (version 19)³⁹. For RNA-seq and Ribo-seq, a maximum of 4 mismatches were
17 allowed and multi-mapping to up to 8 different position was permitted. Alignments flagged as
18 secondary alignment were filtered out, ensuring 1 genomic position per aligned read. Ultimately, we
19 thus obtained ~25 mio aligned reads. To infer the P-site locations, we created a histogram of distances
20 between the 5' end of Ribo-seq reads and well-annotated start and stop codons (CCDS, see below), for
21 each read length (Fig. 1a, Supplementary Fig. 1). Read lengths and offsets used to infer the P-sites
22 position are available in the Supplementary Table 1 for the different Ribo-seq libraries used. RNA-sites
23 were calculated using an arbitrary (26th) position for each RNA-seq read. TPM values for RNA-seq and
24 Ribo-seq were calculated using RSEM⁴⁰. Custom Unix scripts and BedTools⁴¹ were used to create data
25 tracks containing 1) Ribo-seq coverage; 2) P-sites distributions, 3) RNA-seq coverage and 4) RNA-sites
26 distributions for different genomic regions.

27

28 **Exon-level annotation, simulations and ORF identification**

29 Data tracks were created for each annotated exon in the GENCODE v19 annotation, distinguishing
30 between regions annotated as consensus coding sequences (CCDS), non-CCDS exons inside CCDS-
31 containing genes, and exons in non-CCDS genes. Non-CCDS regions (5'UTRs, alternative exons etc...)
32 were annotated with respect to CCDS locations. Exons with more than 5 P-sites were considered for

1 quality control checks. For the benchmarking test, we sampled 1000 CCDS exons from different read
2 lengths and coverage as a positive set. For each exon, we randomly shuffled 1000 times the P-sites
3 positions to obtain a negative set.

4

5 Full ORFs were detected by merging exons according to the transcript structures reported in the
6 GENCODE v19 annotation. For CCDS genes, all CCDS transcripts comprising the “appris” tag were used,
7 by prioritizing transcript with the “appris_principal” tag. For non-CCDS transcripts, all annotated
8 transcripts containing an exon with >5 P-sites were used. For Zebrafish, all transcripts structures
9 annotated in Ensembl (version 76) were used.

10

11 For each transcript, every pair of consecutive start-stop codons (ORF) was tested for its 3nt periodic
12 pattern using the multitaper method, in all the three possible frames (p-value <0.05). ORFs with less
13 than 50% of in-frame P-sites were excluded. In case of multiple possible start codons, we chose the most
14 upstream in-frame ATG with more than 5 P-sites positions (>50% in-frame) between it and its closest
15 neighbor ATG (Fig. 2a). In case of multiple transcript isoforms harboring the same ORF, the transcript
16 with the best support from RNA-seq was chosen.

17

18 ORFs were annotated as follows: ORFs_ccds as ORFs overlapping known CDS regions in CCDS genes;
19 non-CCDS coding ORFs were defined similarly, but when overlapping non-CCDS CDS regions. uORFs and
20 dORFs as ORFs in CCDS genes non overlapping with any CDS exon, and annotated with respect to the
21 annotated transcript CDS; ncORFs as ORFs in non-CCDS genes and not overlapping with any CDS exon
22 (Supplementary Fig. 7). ORFs were filtered out when >30% of the Ribo-seq coverage supported by multi-
23 mapping reads only (filtering was disabled for the custom peptide database creation).

24

25 **Multitaper method**

26 In digital signal processing, the quantitative estimation of periodic components in a finite signal (the
27 power spectral density, or PSD) is an intense area of study, with application to diverse fields of scientific
28 research. A switch from the original representation of a signal (the time domain) to its spectrum of fixed
29 periodic components (the frequency domain) is achieved via the Fourier Transform. In the frequency
30 domain, a vector of coefficients represents the contribution of each frequency component in shaping
31 the original signal.

32

1 The raw output of the Fourier transform (the periodogram) typically suffers from high variability, and as
2 such, it represents a poor estimate of the PSD. Moreover, the limited amount of available realizations of
3 the same signal (i.e. the lack of replicates) poses a real challenge in the estimation of robust coefficients
4 for different frequencies. Applying a smoothing window (taper) to the signal before calculating its
5 Fourier transform helps reducing such variability, but generally creates a biased estimate of the PSD. The
6 choice of taper functions is fundamental, and many different solutions have been proposed in the last
7 decades to find the “optimal” window function.

8

9 The multitaper method, originally proposed by Thomson¹⁸, offers a promising, non-parametric solution
10 to the PSD estimation problem. Its central idea is to apply a set of multiple window functions to the
11 signal, and average the spectral estimates of the ensemble of tapered signals. As the window functions
12 used in the multitaper method are a set of orthonormal functions, the resulting spectra are independent
13 and can be averaged. Specifically, discrete prolate spheroidal sequences (dpss, or Slepian sequences)
14 have been shown to maximize the information content of a finite signal at a given frequency resolution.
15 The use of the Slepian sequences in the multitaper method allows for an optimal solution for reducing
16 the variance in the power spectrum.

17

18 The use of multiple orthonormal window functions on the same signal also enables us to test the
19 robustness of the estimated spectral coefficients. The amount of variance captured by the estimated
20 coefficient at each frequency bin can be compared against a null hypothesis of white noise, leading to a
21 reliable statistical test to determine significant frequency components²⁰.

22

23 **RiboTaper**

24 The original multitaper algorithm from Thomson is implemented in R in the package “multitaper”⁴². A
25 stretch of zeros was added to the input sample to reach a minimum length of 50 nt. The multitaper was
26 run with 24 tapers, setting the time-window parameter to 12. Moreover, sequences shorter than 500 nt
27 were zero-padded to 1024 data points before computing the Discrete Fourier Transform, to obtain an
28 adequate frequency resolution in the spectrum. F-values were extracted from the frequency bin closest
29 to 3nt periodicity. P-values from the F-statistic²⁰ were calculated by using 2 and 2k-2 degrees of
30 freedom, where k is the number of tapers (24 in this study). ORFs and exons with less than 6 P-sites or
31 shorter than 6 nt were ignored.

32

1 QTI-seq comparison

2 For every reported QTI-seq peak²², we selected the closest ORF called by RiboTaper based on the
3 reported distance relative to the annotated start codon. Only ATG start codons were used.

4

5 Conservation Analysis

6 PhastCons scores were extracted as average over the entire ORF, or in 25nt windows around start and
7 stop codon. ORFs were then scored with PhyloCSF in the "mle" (default) mode, using the "29mammals"
8 parameter set on the 46-vertebrate alignment to the human genome (hg19) after alignment filtering
9 steps as described in Bazzini *et al*⁶. We additionally used the hexamer score from the CPAT tool to
10 assess the coding potential of different ORFs, using the available trained model for the human genome.
11 For each category, the scores were compared against a control set of ORFs matching length and
12 conservation of the category of interest. For ORFs_ccds and nonccds coding ORFs, we selected ORFs
13 shorter than 300nt as meaningful matching controls (Supplementary Fig. 10). SNPs were downloaded as
14 .gvf files from Ensembl (v75, 1000 Genomes phase 1). We removed SNPs in reverse orientation, SNPs
15 falling into genomic repeats (using the RepeatMasker track from the UCSC genome browser, March
16 22, 2015), and rare SNPs with derived allele frequency <1%. We then recorded for each ORF and its
17 conceptual translation the number of synonymous and nonsynonymous SNPs when comparing to the
18 human reference genome, as well as the number of synonymous and nonsynonymous sites derived
19 from the degeneracy of the genetic code. For every set of ORFs, we aggregated these numbers and
20 calculated the dN/dS ratio, where dN is the number of nonsynonymous SNPs per nonsynonymous site,
21 and dS the number of synonymous SNPs per synonymous site, respectively. For the CPAT and PhyloCSF
22 scores, p-values come from Wilcoxon-Mann-Whitney tests. For dN/dS ratios, p-values come from a Chi-
23 square test, using as expected frequencies the values from the ORF control set⁴³.

24

25 Mass-spec sample collection and preparation

26 The proteomic data for HEK293 was published recently²⁷ (PRIDE accession number: PXD002389). Briefly,
27 cells were grown in DMEM (life technologies). Lysis was performed in 50 mM ammonium bicarbonate
28 buffer (ABC, pH 8.0) containing 2% SDS and 0.1 M DTT. Sulfhydryl groups were alkylated by adding
29 iodoacetamide to a final concentration of 0.25 M and incubation for 20 min. Proteins were precipitated,
30 resuspended in 6M urea/ 2M thiourea/ 10mM HEPES and digested into peptides using LysC (3 h) and
31 Trypsin (overnight, diluted 4x with 50 mM ABC). Peptides were then acidified, desalted and subjected to
32 isoelectric focusing (IEF) for fractionation.

1

2 **LC-MS/MS and data analysis**

3 Peptides were desalted using stage tip purification and subsequently analyzed by online liquid-
4 chromatography tandem mass-spectrometry on a Q-Exactive (ThermoFisher) instrument using nano-
5 electrospray ionization. Resolution was set to 70,000 and 17,500 for full and fragments scans
6 respectively. Peptides were identified from MS/MS spectra by searching against the recent Uniprot
7 human database (2014-10) or the newly generated HEK293 specific database using ribosome profiling
8 using MaxQuant⁴⁴ version 1.5.2.8. For all searches carbamidomethyl (C) was set as fixed, oxidation (M)
9 and acetylation (protein N-term) and deamidation (NQ) as variable modifications. A maximum of two
10 missed cleavages was allowed. Peptide FDR was set to 0.01, minimal peptide length was set to 7 amino
11 acids and the main search peptide tolerance was set to 4.5 ppm. The protein FDR was disabled.

12

13 **Mass-spec data processing**

14 Custom peptide databases were built by using all the set of identified ORFs prior to filtering for multi-
15 mapping reads. FDR was calculated based on the ratio of hits in the positive and decoy database, as
16 previously described⁴⁴. Counts and feature distribution (PEP, Score, Sequence length) from evidence files
17 were compared based on a FDR < 1 %, excluding Reverse Hits and Contaminants as well as using unique
18 sequence information. Non-Uniprot peptide sequences were defined using PeptideMatch⁴⁵.

19

20 **Accession codes**

21 Gene Expression Omnibus (GEO) GSE73136 (Ribo-seq data from HEK293 cells).

22

23 **Software availability**

24 All the scripts (written in Unix, Bedtools and R) to create the data tracks, run the multitaper analysis,
25 annotate the exons, perform the simulations, reconstruct the full ORFs etc. are available at
26 <https://ohlerlab.mdc-berlin.de/software> and in the Supplementary Software. The method requires a
27 genome fasta file, a GTF file for annotation, alignment files for Ribo-seq and RNA-seq and the read
28 lengths (with respective cutoffs) used for the P-sites calculation.

29

30

31

32 **Figure Captions:**

1

2 Fig. 1: The RiboTaper workflow.

3 Ribo-seq reads are first mapped to the genome using a split-aware aligner (e.g. STAR). a) To infer P-site
4 positions, aggregate profiles are created over annotated start and stop codons, for different read
5 lengths. Data tracks (P-sites tracks shown) are created for all annotated exons; three examples for a
6 coding exon, UTR exon and exon in a non-coding gene are shown. The tracks are analyzed with the
7 multitaper method, which applies a set of orthonormal functions (known as Slepian sequences, shown
8 for 7 tapers) to allow for robust estimates of every frequency component. Finally a statistical
9 significance test is performed for each exon at a frequency of 3nt. Three examples of length- and
10 coverage-matched exons are derived from different transcript classes and exhibit differences in
11 significant 3 nt periodicity. b) Impact of ORF length and expression on the percent of expressed CCDS
12 exons detected to exhibit significant 3nt periodicity (exact value on each bar). Low, Med, and High bins
13 correspond to 2.34–3.84, 5.91–9.14, and 15.08–29.77 RPKMs respectively. c) Histogram of p-values for
14 all CCDS exons tested applying MultiTaper (top) and Chi-squared (bottom) tests to RNA-seq data as a
15 control that should not exhibit extensive periodicity. P-values < 0.05 are shown in red.

16

17

18 Fig. 2: *De novo* ORF reconstruction and examples of significant ORFs.

19 In case of multiple possible start codons it is necessary to discriminate between internal methionines
20 and initiation codons. a) For a given transcript we anchor the analysis on the stop codon and utilize the
21 most upstream in-frame ATG with more than 5 P-sites positions (>50% in-frame) between it and its
22 closest downstream ATG (green dashed box). Examples of translated ORFs are shown in b) for two
23 uORFs and one CCDS ORF for MIEF1 and in c) for ncORF in CTD-2162K18.5. For a given transcript each
24 plot depicts the RNA-seq coverage (grey), the potential ORFs of the three possible frames (red, green,
25 and blue), and the ORFS with significant PFP colored by the annotation category. All of the PFP ORFs
26 detected are supported by peptides from mass-spec data except for the most 5' uORF in MIEF1.

27

28

29 Fig. 3: Comparative analysis of translated ORFs from coding and non-coding regions.

30 a) Histogram of expression levels for genes with (right) and without (left) translated ORFs by annotation
31 category. b) Number, length and coverage of protein coding ORFs and ncORFs, along the different sub-
32 categories. c) Scatterplot (top) of the distance between reported QTI-seq ATG peaks and annotated start

1 codons (x-axis) vs distance between RiboTaper ORFs starts and the annotation (y-axis). Barplot (bottom)
2 of the number of start positions identified by both QTI-seq and RiboTaper with respect to the annotated
3 translation initiation site (aTIS). d) Overlap (top) and coverage (bottom) of ORFs identified in the Gao *et*
4 *al* data set compared to our data set, split by ORF category.

5

6 Fig. 4: Conserved and non-conserved RiboTaper-identified ORFs define the cellular proteome.

7 a) Local nucleotide conservation (computed by PhastCons) in a +/- 25 nt window around start and stop
8 codons positions. b) Comparison of length- and conservation-matched control ORFs to the different ORF
9 categories, by CPAT hexamer score and c) dN/dS ratio. ORFs_ccds and nonccds coding ORFs <300
10 nucleotides long were used in this analysis, to match negative control ORFs (Supplementary Fig. 10). d)
11 Overlap between the protein databases derived from Uniprot or RiboTaper ORFs for all possible peptide
12 sequences (left) and for detected peptide sequences (right). e) Gene expression levels (Ribo-seq and
13 RNA-seq) for genes showing peptide support in the two search strategies (RiboTaper vs Uniprot). f) The
14 number of genes containing at least one RiboTaper identified ORF with peptide evidence by ORF
15 category (left). Genes containing at least one RiboTaper identified ORF with novel peptide evidence (not
16 found in Uniprot, human entries, rel. October 2014) using the RiboTaper database (middle; 191 ORFs in
17 189 genes) or a database of the union of RiboTaper and Uniprot entries to exclude potential cross-
18 matches (right; 157 novel peptides mapping to 129 ORFs in 127 genes). Red numbers indicate evidence
19 from uniquely assigned peptides.

20

21

22 Supplementary Figure 1: Metagene analysis for different datasets. Aggregate plots in different read
23 lengths (from 25 to 30 nt) are shown, showing distance between 5'ends and annotated start and stop
24 codons. Distinct profiles, in terms of both precision and coverage, emerge in the different datasets. a)
25 HEK293 - This study; b) HEK293, Gao *et al* 2014; c) Zebrafish 5h post-fertilization, Bazzini *et al* 2014.

26

27 Supplementary Figure 2: Quality control steps in the different datasets used. HEK293, this study: a) Ribo-
28 seq and b) RNA-seq coverage of the different used datasets, together with the number of defined P-site
29 positions. c) Number of CCDS exons with more than 5 Ribo-seq and RNA-seq reads. d) % of CCDS exons
30 called by RiboTaper or Chi-squared test using Ribo-seq (p-value lower than 0.05), together with % of
31 negative exons using RNA-seq (p-value higher than 0.05). e) Histogram of p-values for the multitaper

1 and f) Chi-squared test on CCDS exonic RNA profiles. g) Accumulation of P-sites on the maximum frame
2 for exons in different transcripts regions. h) Frame agreement with the CCDS annotation. j) % of periodic
3 exons in different length-coverage categories, derived by taking the second, fourth and sixth of seven
4 quantiles of length-coverage distributions, using Ribo-seq or j) RNA-seq.

5

6 Supplementary Figure 3: Quality control steps in the different datasets used. HEK293, Gao et al: a) Ribo-
7 seq and b) RNA-seq coverage of the different used datasets, together with the number of defined P-site
8 positions. c) Number of CCDS exons with more than 5 Ribo-seq and RNA-seq reads. d) % of CCDS exons
9 called by RiboTaper or Chi-squared test using Ribo-seq (p-value lower than 0.05), together with % of
10 negative exons using RNA-seq (p-value higher than 0.05). e) Histogram of p-values for the multitaper
11 and f) Chi-squared test on CCDS exonic RNA profiles. g) Accumulation of P-sites on the maximum frame
12 for exons in different transcripts regions. h) Frame agreement with the CCDS annotation. j) % of periodic
13 exons in different length-coverage categories, derived by taking the second, fourth and sixth of seven
14 quantiles of length-coverage distributions, using Ribo-seq or j) RNA-seq.

15

16 Supplementary Figure 4: Quality control steps in the different datasets used. Zebrafish, Bazzini et al: a)
17 Ribo-seq and b) RNA-seq coverage of the different used datasets, together with the number of defined
18 P-site positions. c) Number of CDS exons with more than 5 Ribo-seq and RNA-seq reads. d) % of CDS
19 exons called by RiboTaper or Chi-squared test using Ribo-seq (p-value lower than 0.05), together with %
20 of negative exons using RNA-seq (p-value higher than 0.05). e) Histogram of p-values for the multitaper
21 and f) Chi-squared test on CDS exonic RNA profiles. g) Supplementary Figure 6: Accumulation of P-sites
22 on the maximum frame for exons in different transcripts regions. h) Frame agreement with the CDS
23 annotation. j) % of periodic exons in different length-coverage categories, derived by taking the second,
24 fourth and sixth of seven quantiles of length-coverage distributions, using Ribo-seq or j) RNA-seq.

25

26

27 Supplementary Figure 5: Multitaper performance for different numbers of tapers. Shown are the a)
28 AUC, b) sensitivity and c) specificity values for real and shuffled exonic P-sites profiles, along different
29 length-coverage categories, derived by taking the second, fourth and sixth of seven quantiles of length-
30 coverage distributions. Little change is detectable by using more than 24 tapers.

1

2

3 Supplementary Figure 6: Different methods performances on P-sites and RNA-sites tracks for coding
4 exons of different length/coverage values. Shown are different ORFscore cutoffs, Chi-squared and
5 multitaper statistical tests (see text for more details). a) Data in HEK 293 cells from our study (CCDS
6 annotation); b) Data in Zebrafish 5hpf from Bazzini et al, 2014.

7

8

9 Supplementary Figure 7: Schematics of RiboTaper ORF annotation. In a) a "uORF" is defined as upstream
10 of the annotated start codon and non-overlapping any coding exon, while different "ORFs_ccds" are
11 overlapping annotated coding exons. A "dORF" is defined as downstream of the stop codon and not
12 overlapping any coding exon. Shown is also a lincRNA ORF overlapping a coding exon, therefore
13 annotated as "nonccds_coding_ORF". In b) a "nonccds_coding_ORF" in a non-CCDS protein coding gene,
14 defined as overlapping a coding exon. A "nonccds_coding_ORF" in a processed_transcript gene is also
15 present. An ncORF is defined as an ORF in a non-CCDS gene not overlapping any coding exon, here in an
16 antisense gene. In c) an ncORF in a lincRNA gene is shown.

17

18 Supplementary Figure 8: FLOSS scores for ORFs identified by RiboTaper. Shown are FLOSS scores with
19 their cumulative distributions for a) CCDS genes and ORFs_ccds. b) 5'UTRs and uORFs, c) 3'UTRs and
20 dORFs, d) non-coding genes and different ncORFs categories. FLOSS values and cutoffs were calculated
21 as in Ingolia et al, 2014. Low FLOSS scores indicate a protein-coding like fragment length distribution.

22

23 Supplementary Figure 9: The toddler ncORF. Shown are P-sites position, RNA-seq coverage and ORF
24 position. Data from Bazzini et al, 2014.

25

26 Supplementary Figure 10: Length and conservation-matched ORFs for the conservation analysis. No
27 significant difference in terms of a) length and b) nucleotide conservation (PhastCons) was found
28 between detected ORFs and randomly chosen controls ORFs. c) Scores from PhyloCSF (used in the -mle
29 mode) agree with the ones from CPAT (see main text and Figure 4 for more details). * = p-value<0.05;

1 **= p-value<0.01; ***= p-value<0.001, Wilcoxon-Mann-Whitney test. ORFs_ccds and nonccds coding
2 ORFs were selected if <300 nucleotides long, to match negative controls ORFs.

3

4 Supplementary Figure 11: Additional statistics about Ribotaper- and Uniprot-only identified peptides. a)
5 Overlap between genes with peptide evidence in the different search strategies. b) Overlap of Peptide
6 Spectrum Matches (PSMs) in the two strategies. c) FDR vs. PSMs count for the two search strategies. d)
7 Comparisons of RiboTaper-only identified peptides vs. Uniprot-only identified peptides (PEP=Posterior
8 Error Probability).

9

10 Supplementary Figure 12: Genomic locations of “non-canonical” ORFs with peptide evidence. A dORF in
11 the SMIM20 gene. Below a screenshot from the GWIPS-viz Genome Browser.

12

13 Supplementary Figure 13: Genomic locations of “non-canonical” ORFs with peptide evidence. A lincRNA
14 ncORF in the LOC93622 gene. Below a screenshot from the GWIPS-viz Genome Browser. A CUG-start
15 codon is supported by initiating ribosomes from QTI-seq and LTM data. The second exon (on the right)
16 showing a strong accumulation of RNA-seq reads is present in circBase, <http://www.circbase.org/>,
17 (hsa_circ_0069092).

18

19 Supplementary Figure 14: Genomic locations of “non-canonical” ORFs with peptide evidence. An
20 antisense ncORF in the CTD-2162K18.5 gene, in a non-conserved genomic region. Below a screenshot
21 from the GWIPS-viz Genome Browser.

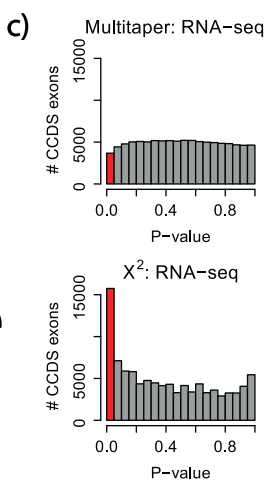
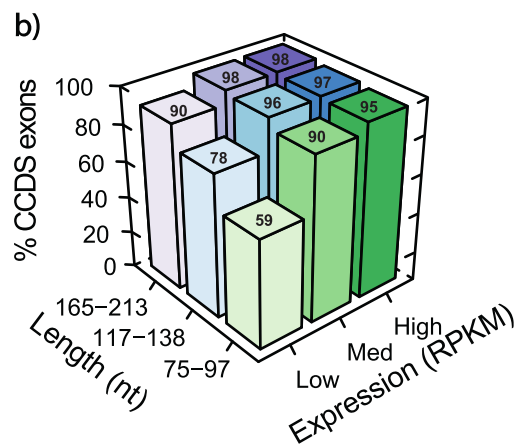
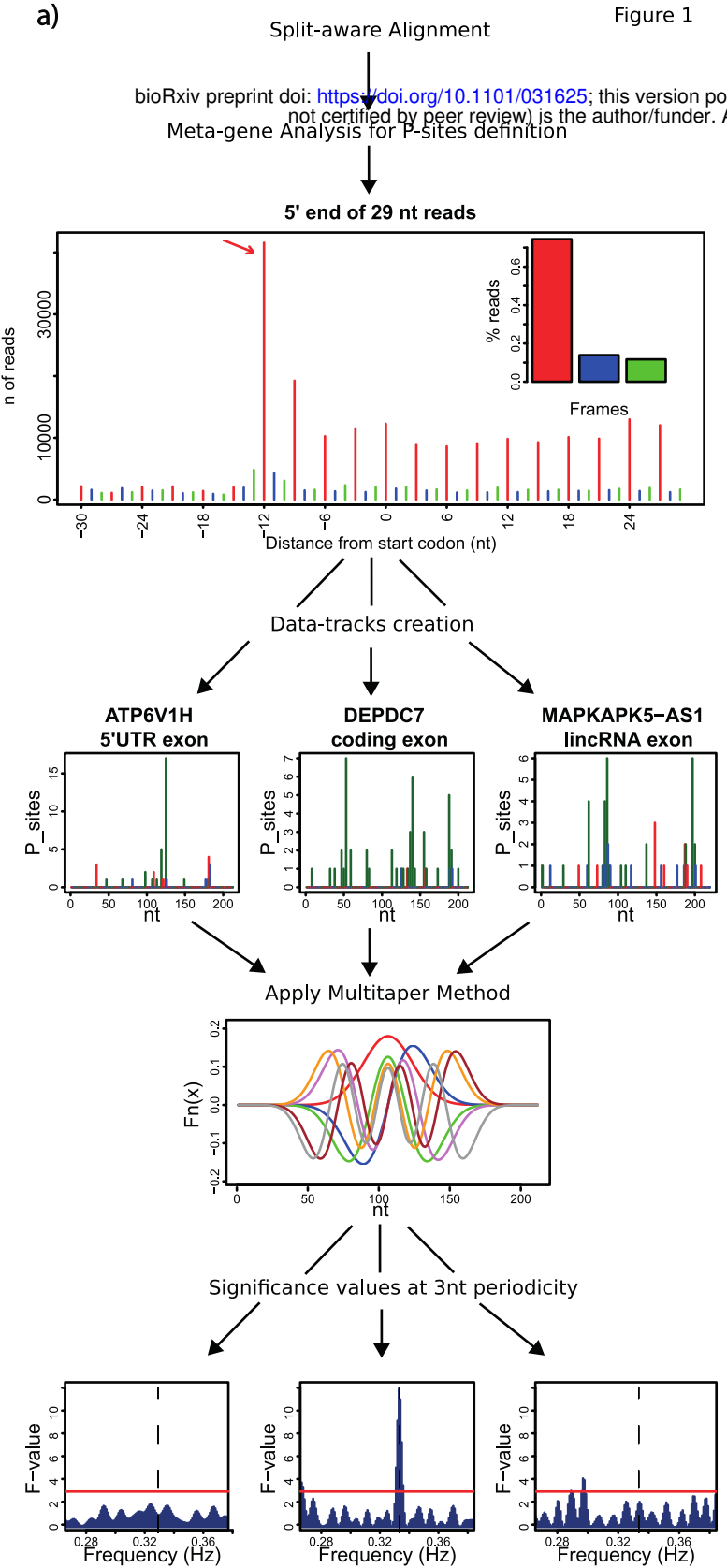
22

23 Supplementary Figure 15: Genomic locations of “non-canonical” ORFs with peptide evidence. An
24 antisense ncORF in the RP11-139H15.1 gene (on the left), in a conserved genomic region. Below a
25 screenshot from the GWIPS-viz Genome Browser.

26

27 Supplementary Table 1: Statistics about alignment and pre-processing of the sequencing datasets used
28 in this study.

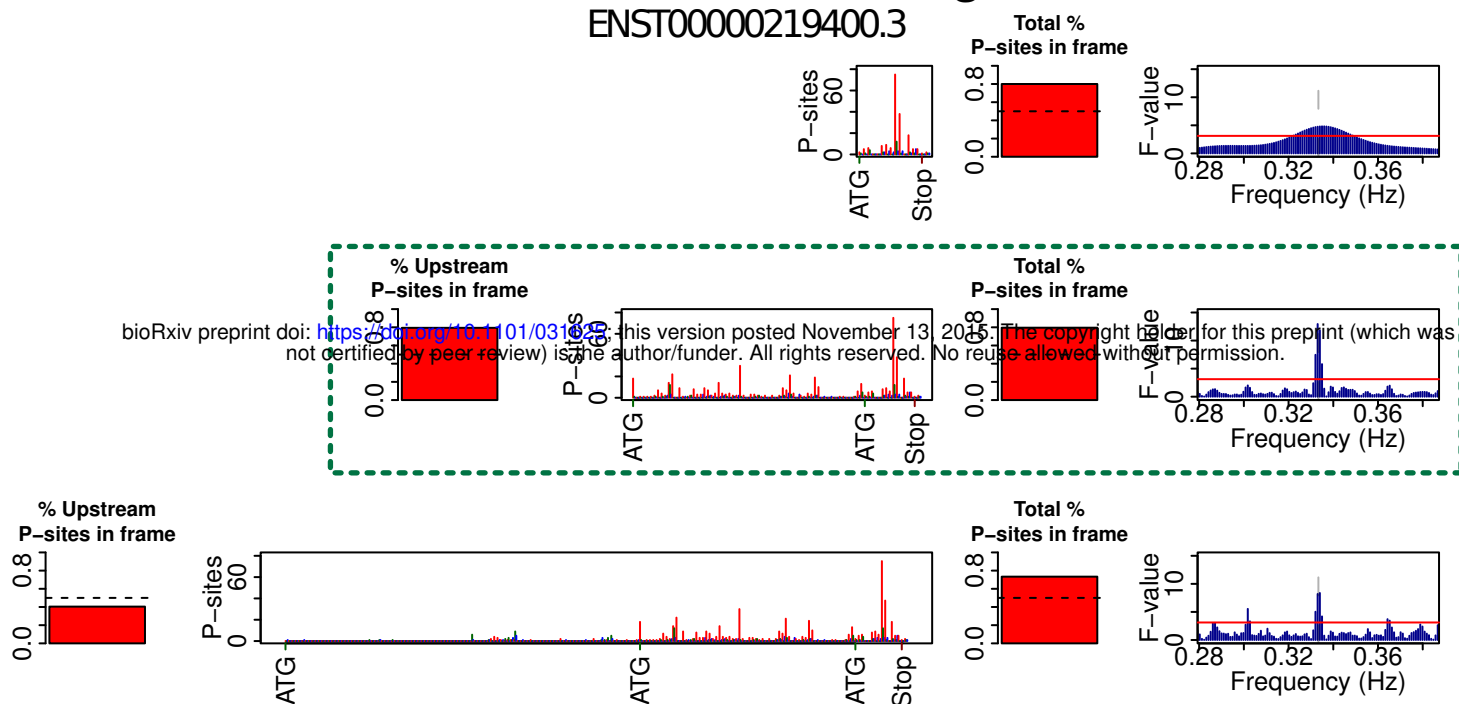
- 1
- 2 Supplementary Table 2: Complete list of identified ORFs in the different datasets used.
- 3
- 4 Supplementary Table 3: Filtered list of identified ORFs in the different datasets used. Non-coding ORFs
- 5 overlapping CDS regions were removed. Additionally, ORFs were filtered out when >30% of the Ribo-seq
- 6 coverage were supported by multi-mapping reads only.
- 7
- 8 Supplementary Table 4: Evidence table for peptides identified by the RiboTaper database only.
- 9
- 10 Supplementary File 1: Archive containing bed files for the identified ORFs.
- 11
- 12 Supplementary Software: RiboTaper (version 1.0) software code.



a)

de novo ORF finding

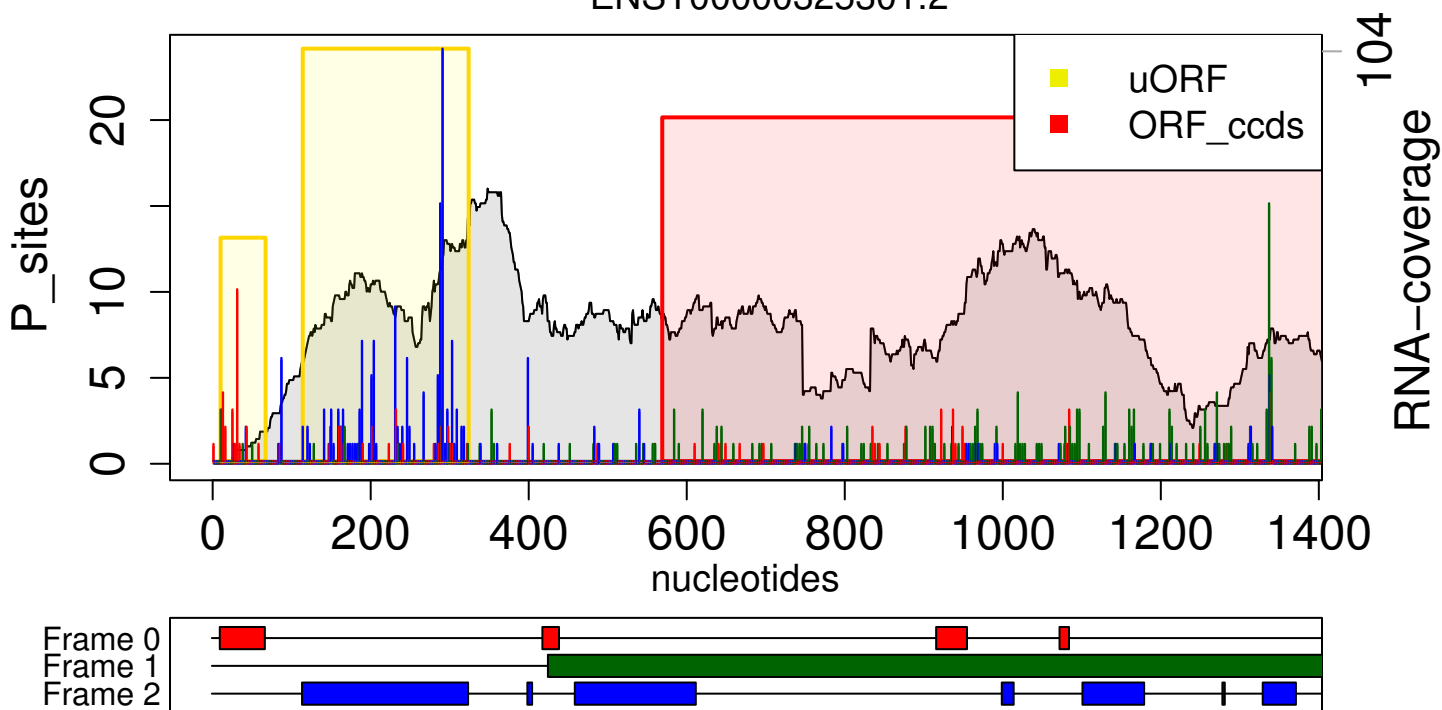
ENST00000219400.3



b)

MIEF1 – protein coding gene

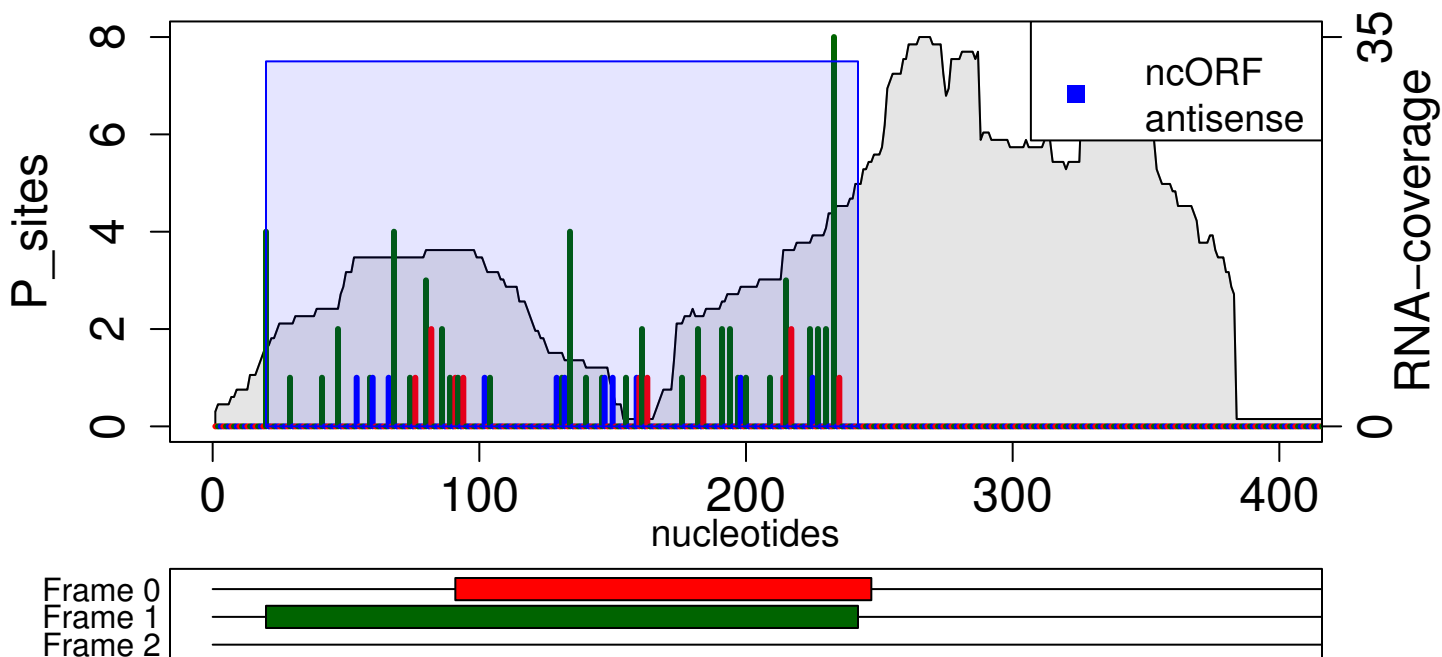
ENST00000325301.2



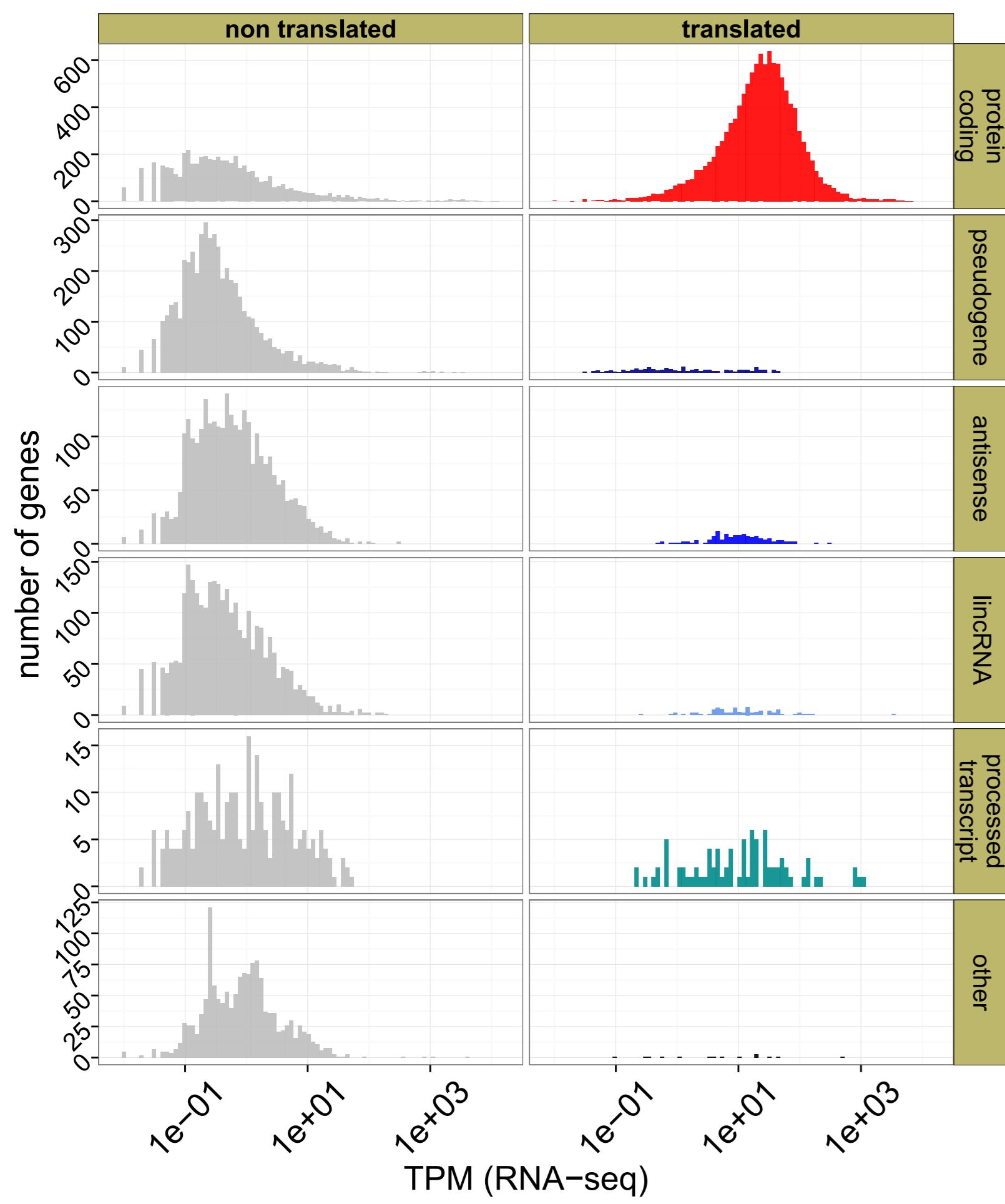
c)

CTD-2162K18.5 – antisense gene

ENST00000589556.1

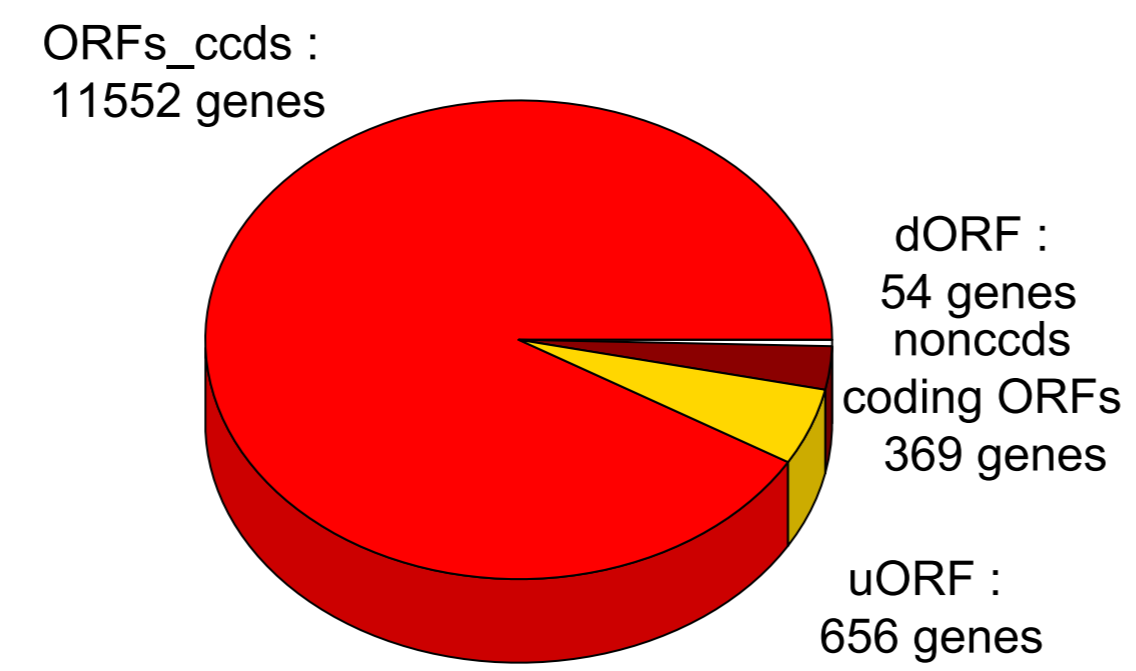


a)

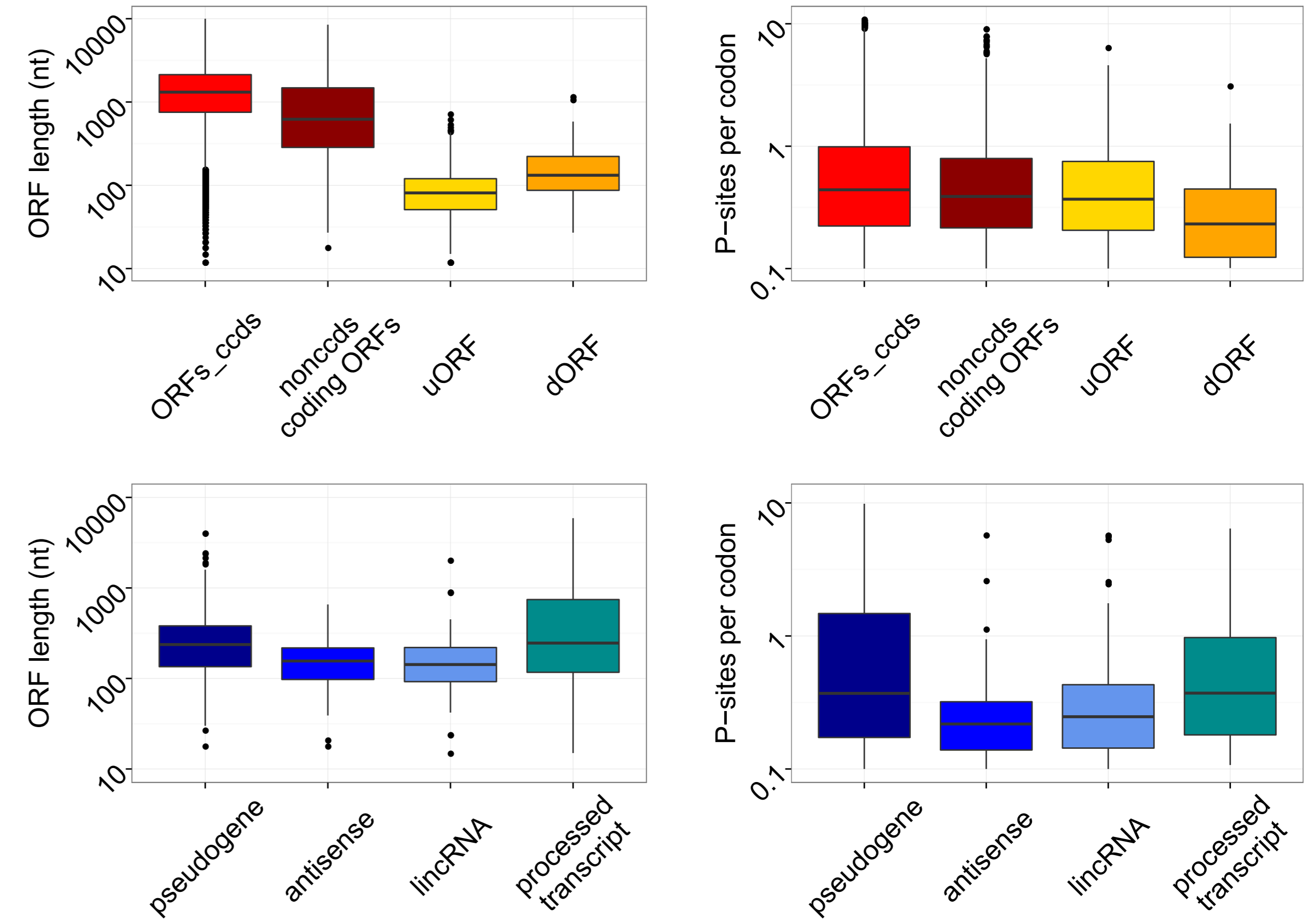
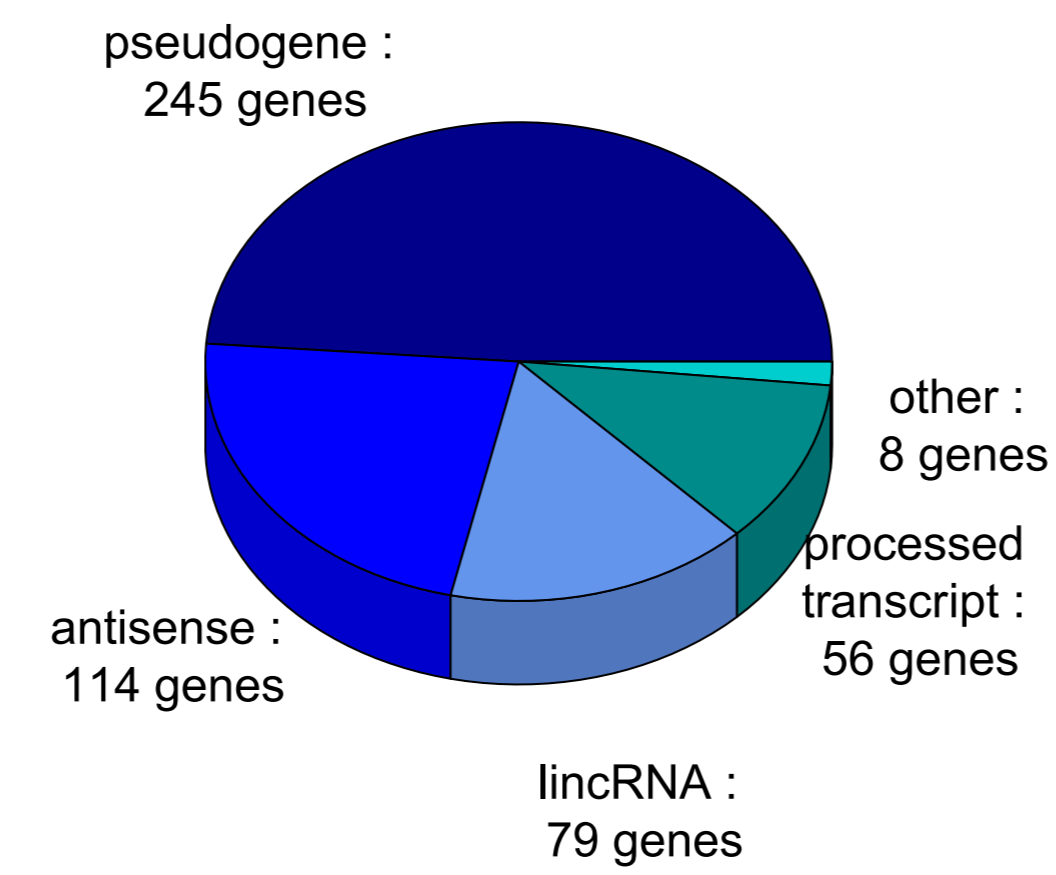


b)

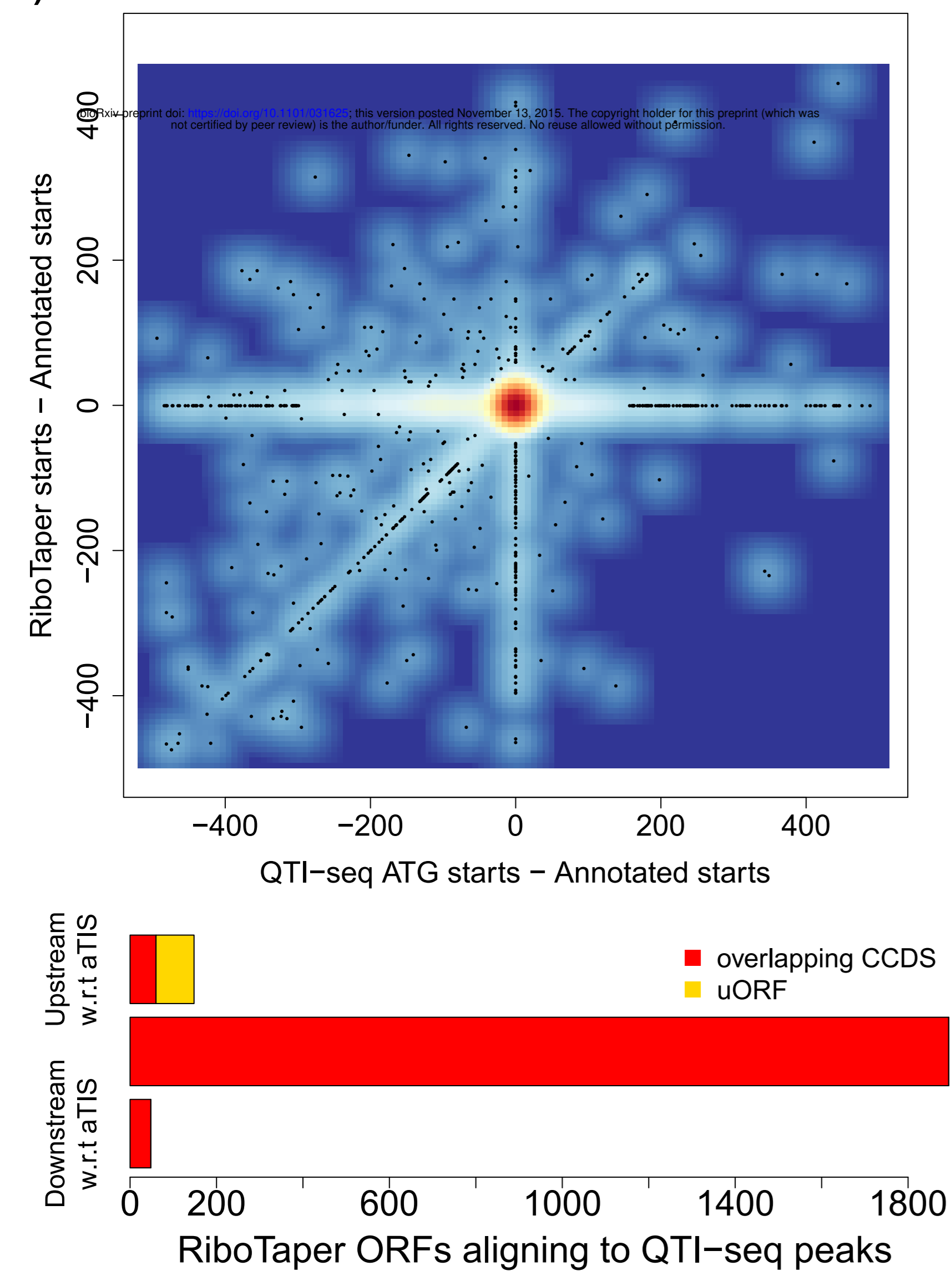
protein coding ORFs



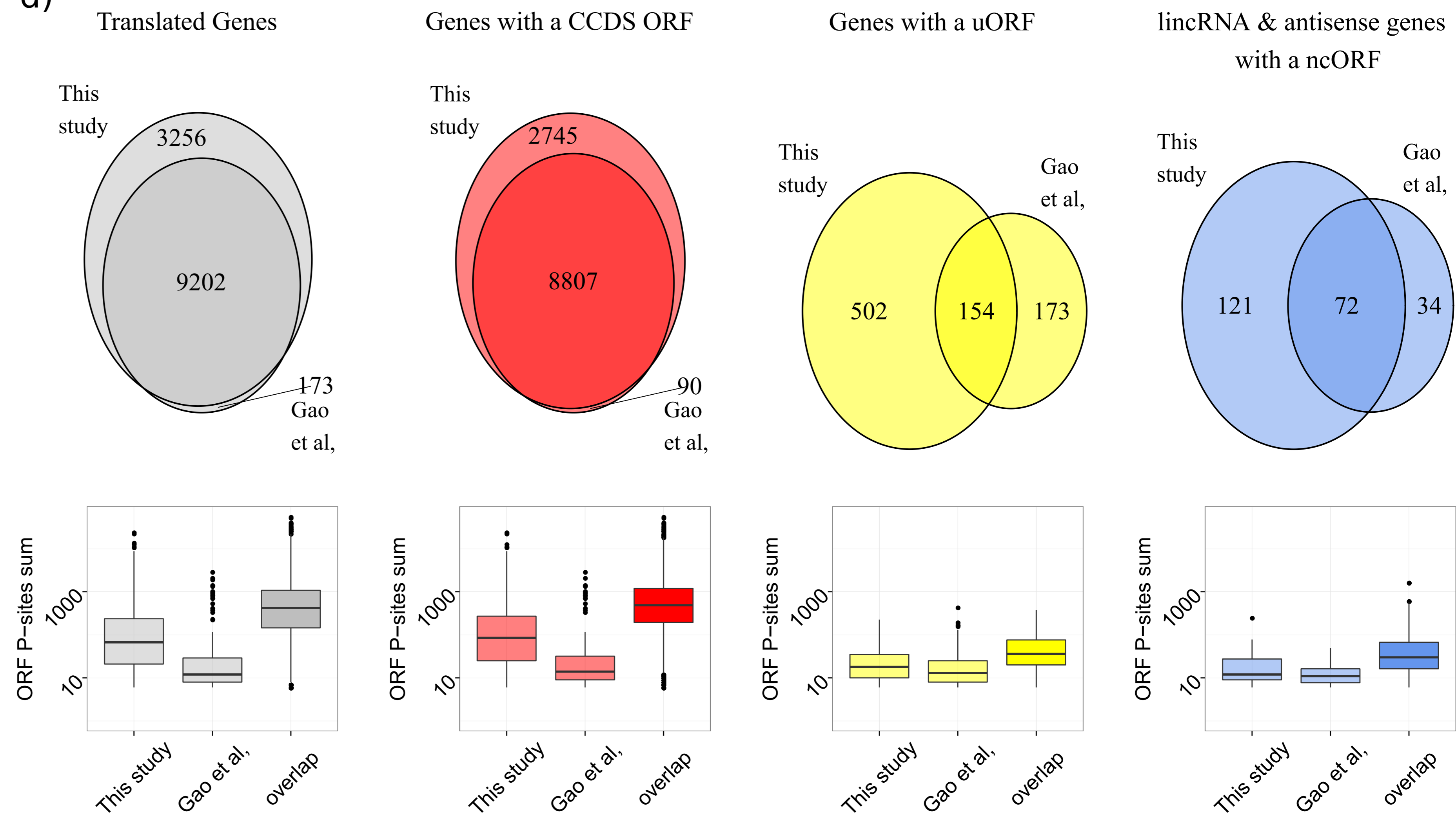
ncORFs

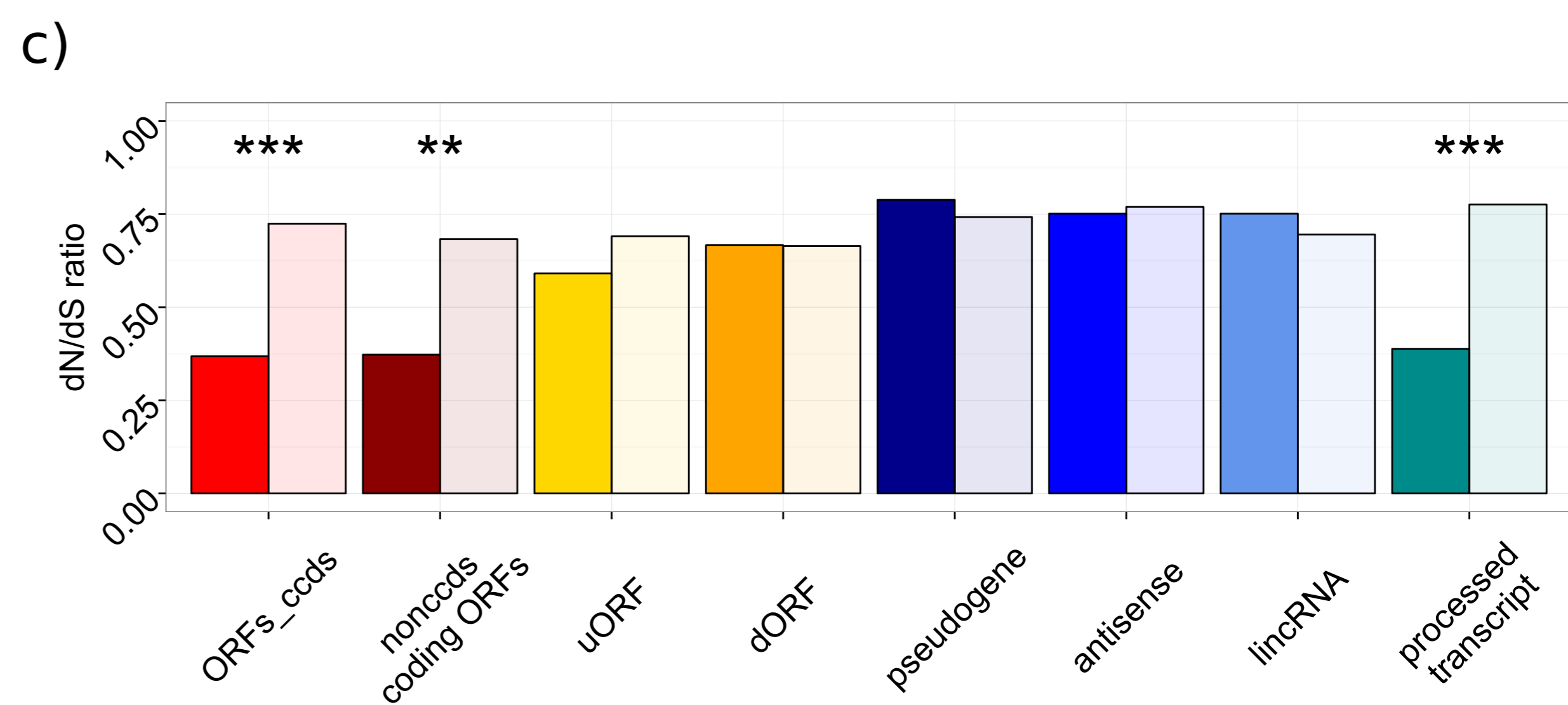
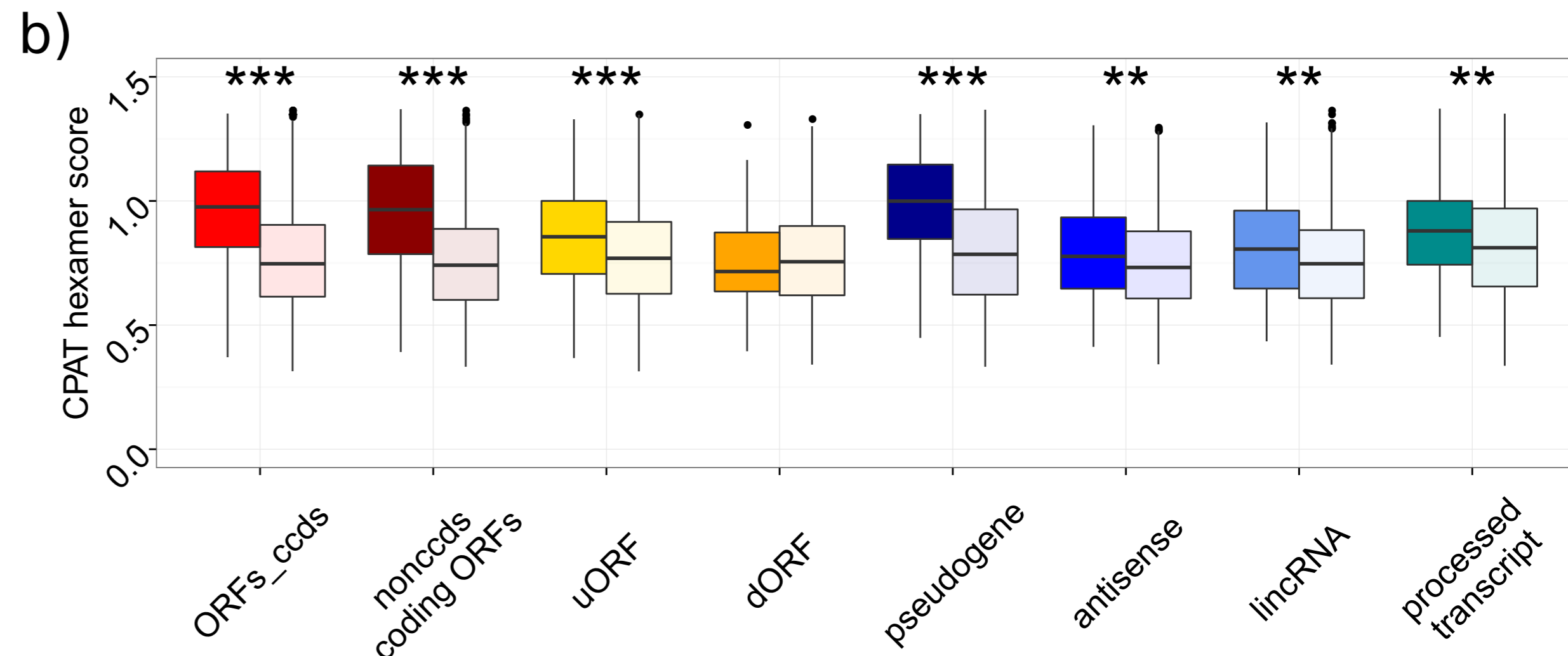
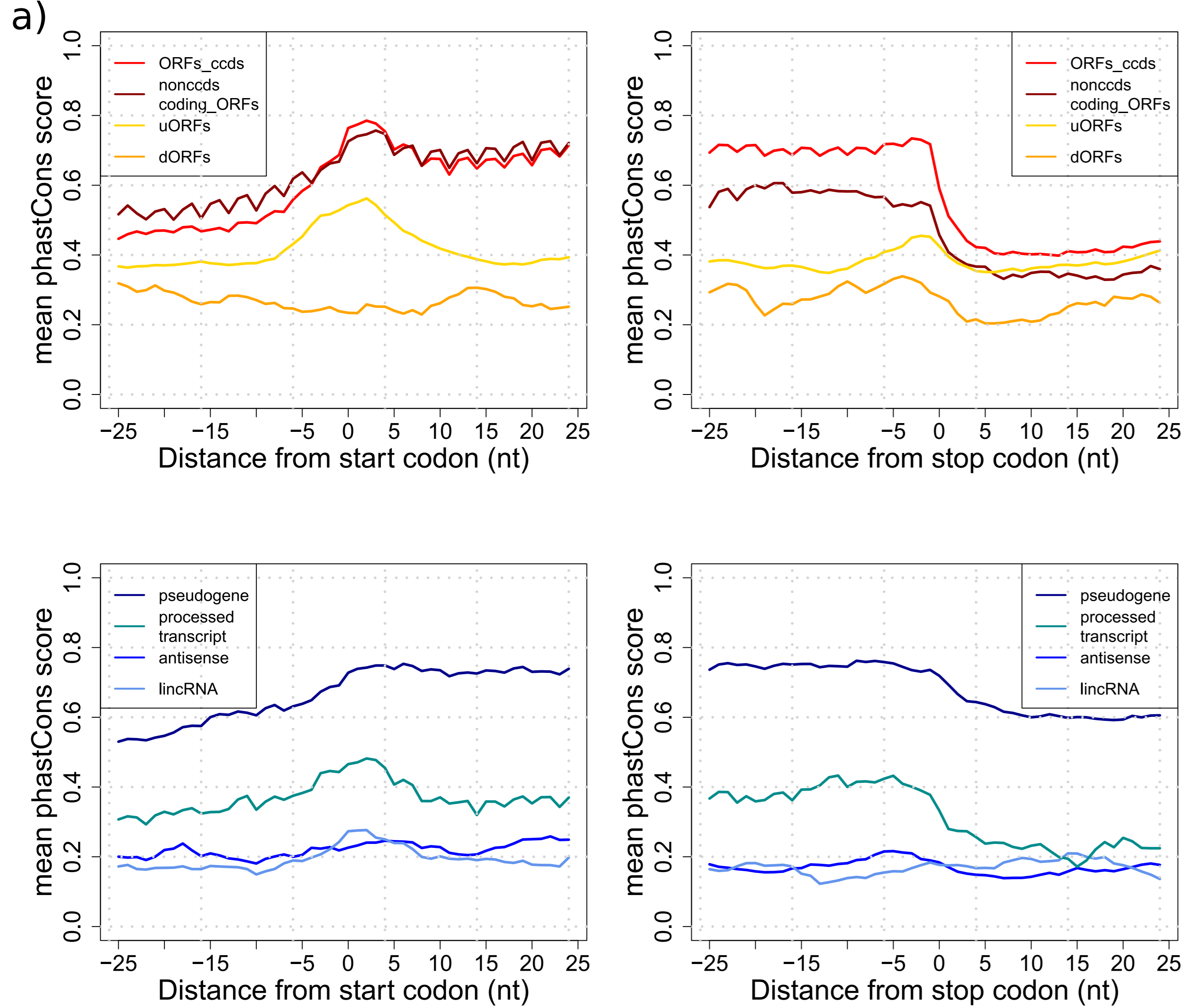


c)

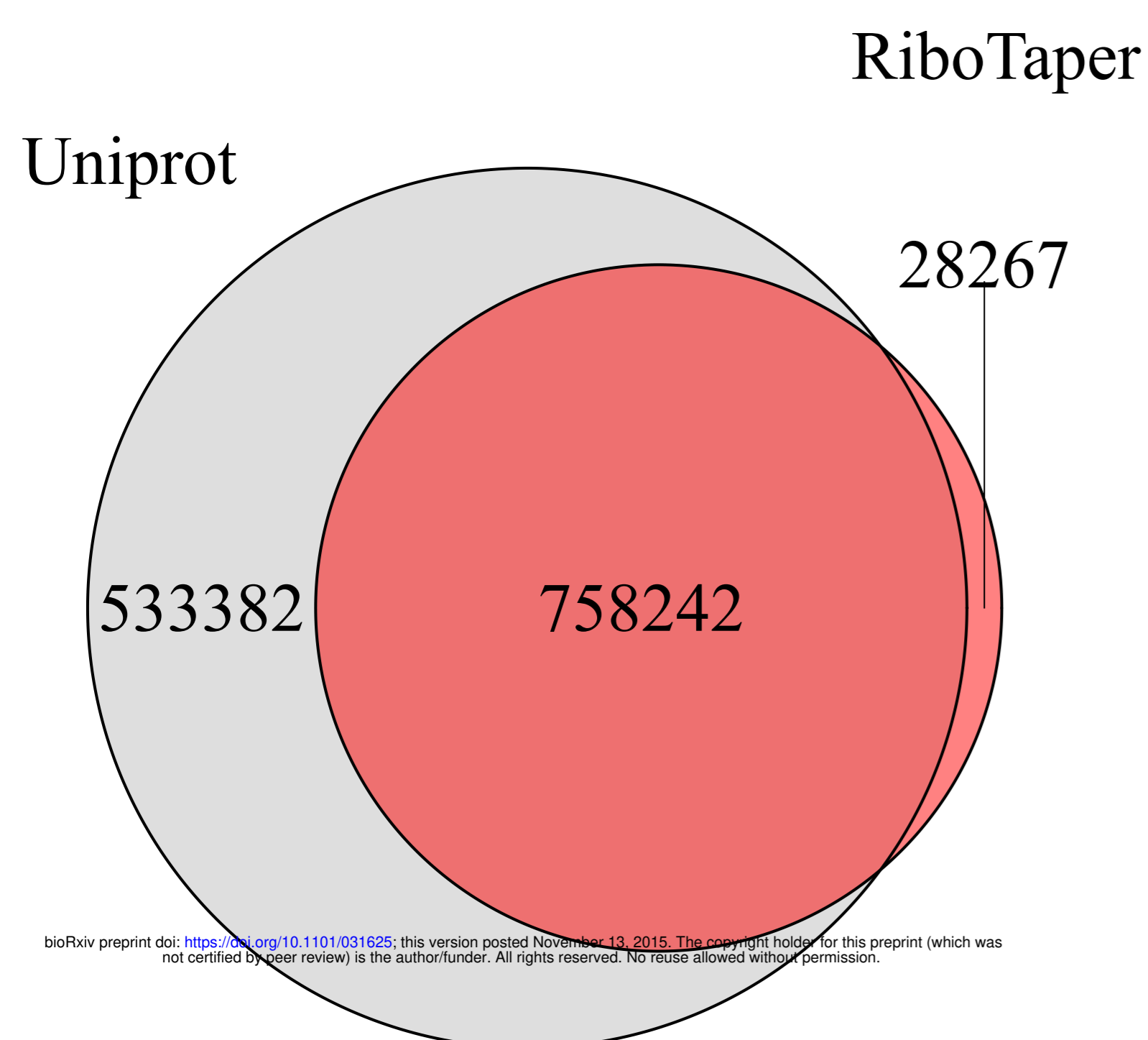


d)

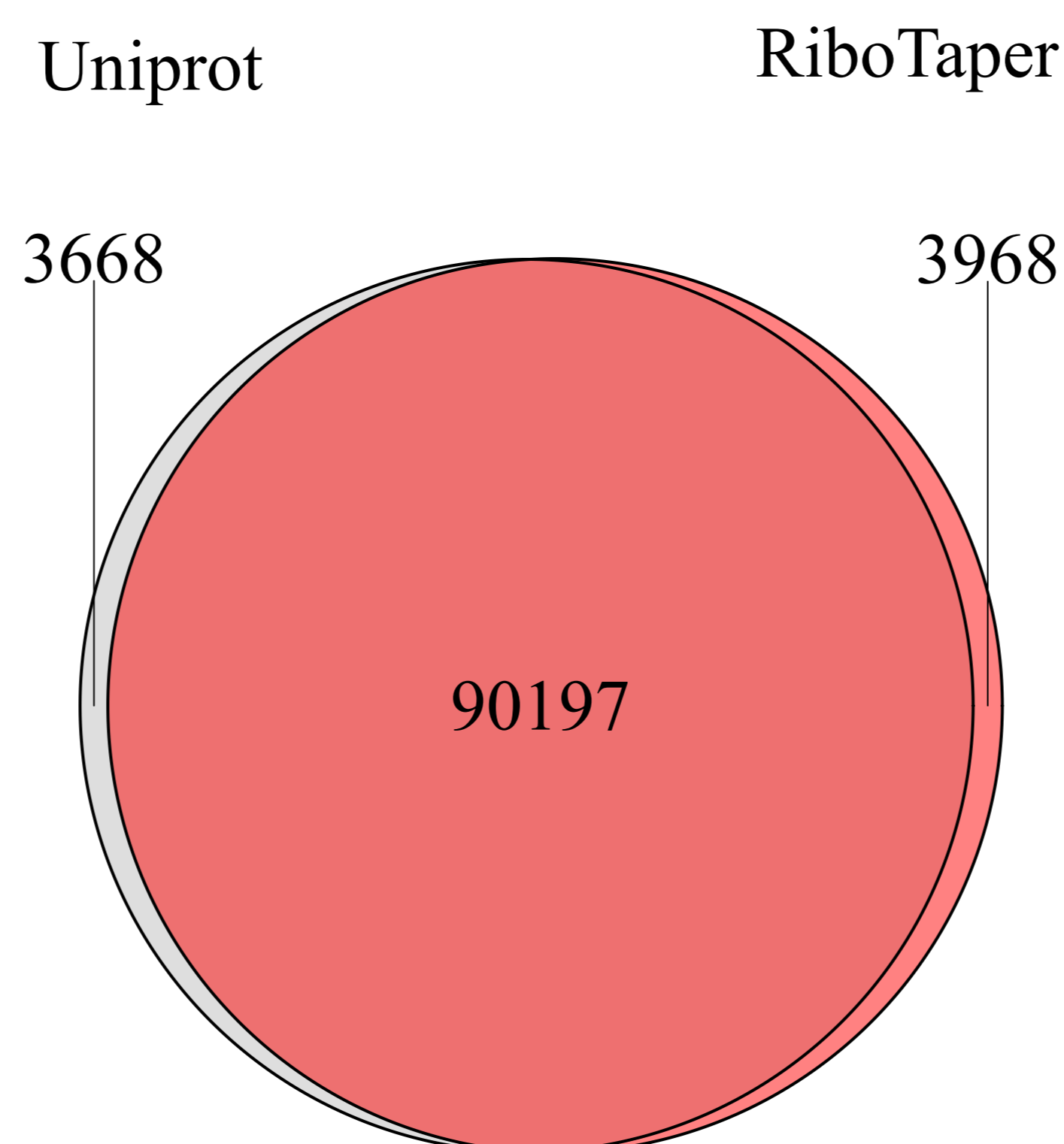




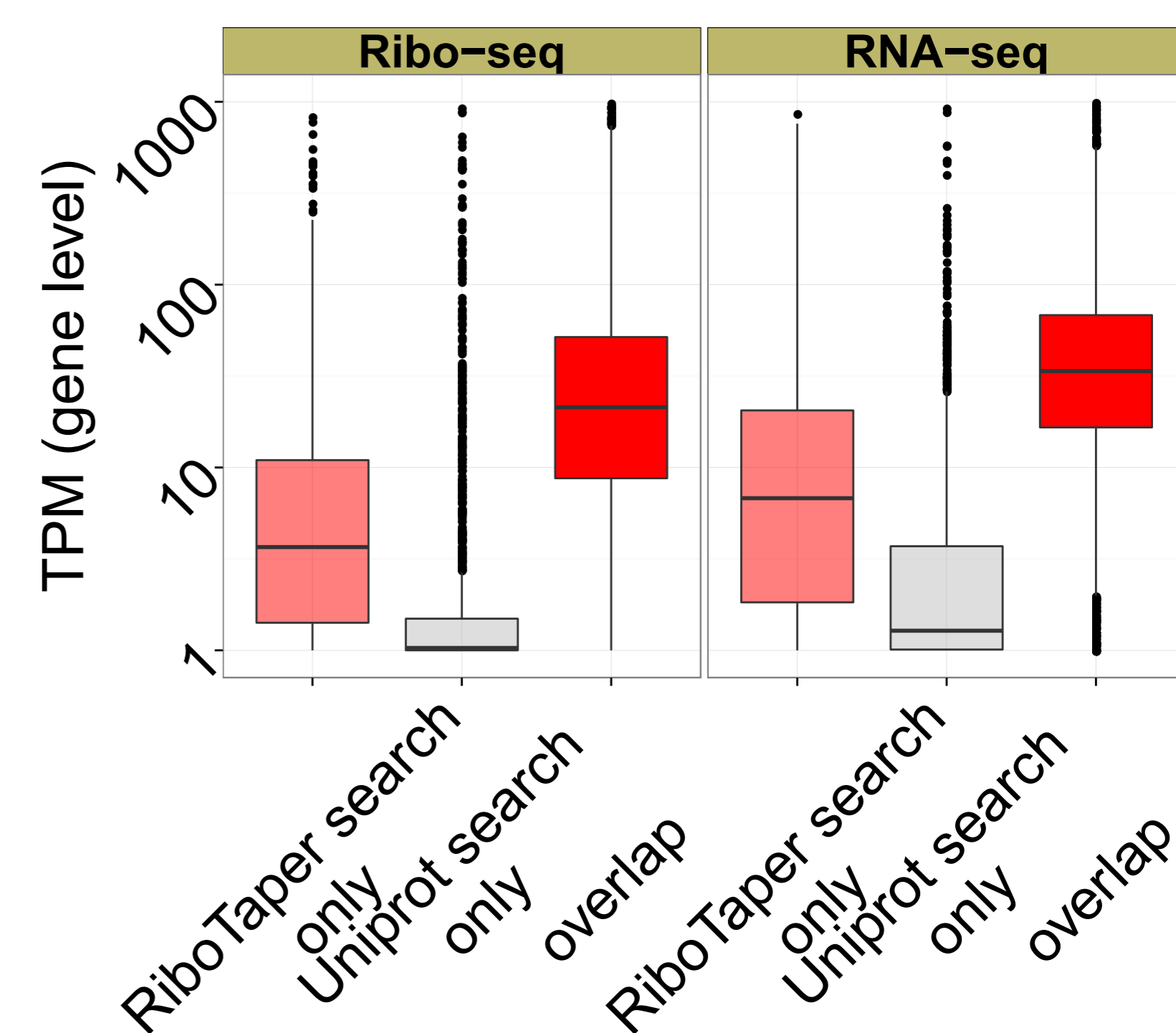
d) in silico digested peptide sequences



Detected peptide sequences



e)



f)

