

Standardized annotation of translated open reading frames

To the Editor — Ribosome profiling (Ribo-seq) has extended our understanding of the translational ‘vocabulary’ of the human genome, uncovering thousands of open reading frames (ORFs) within long noncoding RNAs (lncRNAs) and presumed untranslated regions (UTRs) of protein-coding genes. However, reference gene annotation projects have been circumspect in their incorporation of these ORFs because of uncertainties about their experimental reproducibility and physiological roles. Yet, it is clear that certain ‘Ribo-seq ORFs’ make stable proteins, others mediate gene regulation, and many have medical implications. Ultimately, the absence of standardized ORF annotation has created a circular problem: while Ribo-seq ORFs remain unrecognized by reference annotation databases, this lack of recognition will thwart studies examining their roles. Here, we outline a community-led effort involving Ensembl/GENCODE, the HUGO Gene Nomenclature Committee (HGNC), UniProtKB, HUPO/HPP and PeptideAtlas to produce a standardized catalog of 7,264 human Ribo-seq ORFs; a path to bring protein-level evidence for Ribo-seq ORFs into reference annotation databases; and a roadmap to facilitate research in the global community.

Ribo-seq¹ provides an RNA-sequencing-based readout of mRNA translation by isolating ribosome-bound RNA fragments of ~30 nucleotides in length. Sequencing of these fragments offers genome-wide footprints of ribosome-RNA interactions, detecting translated ORFs with sub-codon resolution^{2–8}. Although Ribo-seq circumnavigates the experimental difficulties of working with protein molecules (for example, using mass spectrometry (MS) analytical tools) and readily finds translations missed by *in silico* evolutionary methods, it does not demonstrate the actual existence of proteins, and most translations do not show signs of constraint as coding sequences (CDS). A wide range of ‘functional’ scenarios are therefore plausible for Ribo-seq ORFs (Table 1).

Several public resources already process and/or display Ribo-seq datasets, including sORFs.org⁹, GWIPS-viz¹⁰ and Trips-Viz¹¹, whereas OpenProt¹² and nORFs.org¹³ incorporate Ribo-seq into

Table 1 | Approaches to interpreting Ribo-seq ORFs

Possible cellular interpretation of Ribo-seq ORF translation	Comments
A Ribo-seq ORF encodes a ‘missing’ conserved protein	Ribo-seq ORFs may be recognized as canonical—in accordance with existing protein annotations—on the basis that the sequence of the proteins they encode shows clear evidence of being maintained by evolutionary selection over a significant period of evolutionary time.
A Ribo-seq ORF encodes a taxonomically restricted protein	Ribo-seq ORFs may encode proteins whose sequences and molecular activities are specific to one species or lineage. Evidence for selection or conservation across distant species or lineages is lacking for these ORFs, either because the protein sequence has diverged beyond recognition from its orthologs, or because the protein evolved recently from previously noncoding material and homologs do not exist in other species or lineages.
A Ribo-seq ORF regulates protein or RNA abundance	Ribosome engagement of regulatory ORFs does not result in a protein product under selection but regulates the abundance of a canonical protein or RNA. This paradigm is well established for uORFs and uORFs, as noted in Table 2, though it is applicable to other transcript scenarios. Regulatory ORFs may compete for ribosomes with their downstream canonical ORFs or produce nascent peptides that stall ribosomes, leading to the controlled ‘dampening’ of protein expression. Alternative modes of action, such as the induction of RNA decay pathways, the processing of small RNA precursors or the adjustment of RNA stability, have also been inferred.
A Ribo-seq ORF is the result of random translation	The translation of some Ribo-seq ORFs may simply be ‘noise’. Because translation has a high bioenergetic cost, a protein that results from random translation is likely to be translated at lower levels than a canonical CDS and evolve neutrally; it may also be comparatively unstable and could be rapidly degraded. Nonetheless, it is theoretically possible that certain proteins do exist as stable ‘junk’ proteins, or that random translation events affect the expression of canonical proteins. The detection of random Ribo-seq ORFs is less likely to be reproducible.
A Ribo-seq ORF encodes a disease-specific protein	This protein would not be produced under normal physiological homeostasis but could be of major interest for diagnostics and therapeutics. Insights of this sort are especially emerging in cancer biology, where transcription and translation are known to be dysregulated. This leads to the production of ‘aberrant’, possibly rapidly degraded proteins that are commonly antigenic and presented on the cell surface by the HLA system, potentially acting as neoantigens. Furthermore, antigens resulting from disease-specific dysregulated ribosome activity—sometimes called defective ribosomal products (DRiPs)—have also been explored.

Note: a given ORF may encompass several of these possibilities: for example, a translated ORF could be both regulatory and implicated in disease neoantigen production.

whole-translatome catalogs. Meanwhile, McGillivray et al. have produced a catalog of upstream ORFs (uORFs) with predicted biological activity¹⁴. Such efforts have made important contributions in Ribo-seq ORF interpretation. Nonetheless, the global scientific community is constrained by the absence of ‘reference’ gene annotation, which supports most large-scale genomics

projects and provides the framework for human variant interpretation (Fig. 1a, Supplementary Fig. 1).

The creation of Ribo-seq annotations within existing reference gene and protein databases presents specific challenges that were not faced by previous cataloging efforts^{9–13}. In particular, it is necessary to consider how these annotations can be

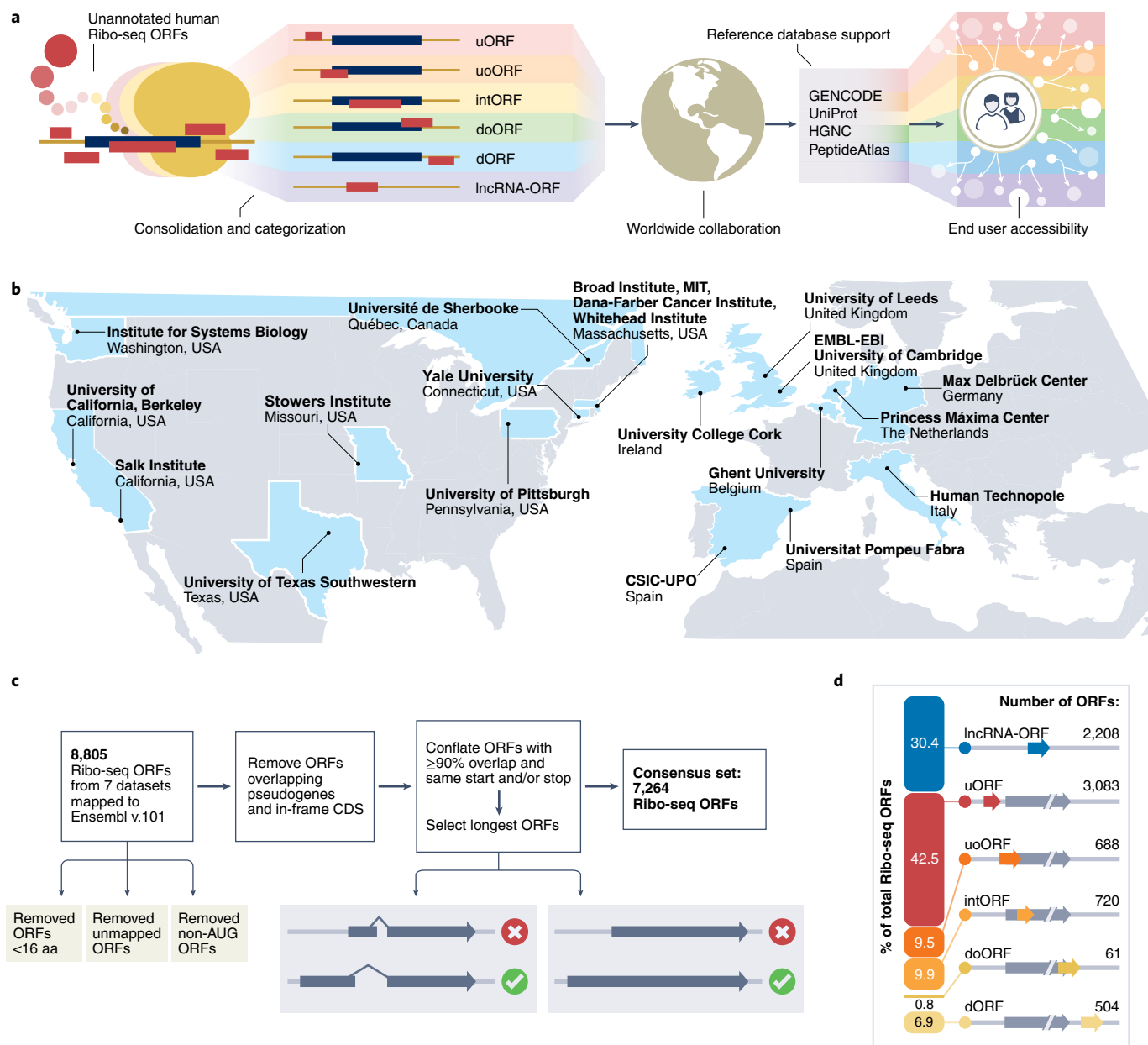


Fig. 1 | Characterization of a consensus set of Ribo-seq ORFs for annotation by GENCODE. **a**, Schematic of the main steps and goals for this consortium effort. **b**, Map showing the participating institutions included in this effort. **c**, Schematic overview of filtering steps used to create the consensus set of ribosome profiling (Ribo-seq) ORFs. **d**, Diagrammatic representation of all Ribo-seq ORFs according to ORF type (see Table 2 for more information).

integrated into the broad range of user workflows that are already supported by global annotation resources. For such reasons, reference annotation projects are generally conservative when it comes to the incorporation of new data types. Thus, rather than attempt to describe a ‘maximal’ set of potential Ribo-seq translations from the outset, our strategy is to build up a comprehensive resource in stages that is reciprocally improved by input from the scientific community (Fig. 1b).

Here, as ‘Phase I’ of this work, we present a consolidated catalog of Ribo-seq ORFs from seven publications^{2–8} annotated onto GENCODE version 35 (Fig. 1c; Supplementary Tables 1–9). A detailed description of the Ribo-seq datasets, our analysis methods and ORF characteristics is available in the Supplementary Methods. We removed ORFs smaller than 16 amino acids (aa) and those translated from non-ATG (‘near-cognate’) initiation codons, and merged redundant sense overlapping

ORFs, resulting in a collated set of 7,264 unique ORFs (Fig. 1c). We classified these ORFs according to their spatial relationship with existing gene annotations (Fig. 1d), as presented in Table 2. We hope community usage of this catalog will help address the key technical and biological questions necessary to move this work into ‘Phase II’, where we aim to create a more comprehensive resource as outlined below.

For Phase I, we investigated repeated ORF identifications between studies,

Table 2 | Terminology and categories of Ribo-seq ORFs

Term	Definition	Biological role(s)
Ribo-seq ORF	Translated sequences identified by the Ribo-seq assay that have not already been annotated by reference annotation projects. Also known as: noncanonical ORFs, alternative ORFs (altORFs), novel ORFs (nORFs) or, when <100 amino acids in size, small ORFs (smORFs), short ORFs (sORFs). Putative encoded proteins in smORFs/sORFs are also known as: microproteins, micropeptides, short ORF-encoded polypeptides (SEPs).	See below.
Upstream ORFs (uORFs)	Translated sequences located within the exons of the 5' untranslated region (UTR) of annotated protein-coding genes.	Regulate translational efficiency of the downstream canonical protein. Cellular-stress-related translation. May produce independently functional proteins.
Upstream overlapping ORFs (uoORFs)	Translated sequences beginning in the 5' UTR of an annotated protein-coding gene and partially overlapping its coding sequence in a different reading frame.	Similar to uORFs. Regulate translation of their overlapping CDS, but with potentially stronger regulatory potential than uORFs. May produce independently functional proteins.
Downstream ORFs (dORFs)	Translated sequences located within the 3' UTR of annotated protein-coding genes	Less commonly detected and generally poorly understood. May act as cis translational regulators.
Downstream overlapping ORFs (doORFs)	Translated sequences beginning in the genomic coordinates of an annotated CDS but continuing beyond the annotated CDS and terminating in the 3' UTR of the annotated protein-coding gene.	Similar to dORFs.
Internal out-of-frame ORFs (intORFs)	Translated sequences located on the mRNA of an annotated protein-coding gene and completely encompassed within the canonical CDS, but translated via a different reading frame. Also known as: altCDSs, nested ORFs, dual-coding regions.	May regulate translation similarly to uORFs in some cases. Detection of intORFs with Ribo-seq is possible but difficult due to the convolution of triplet periodicity signals from two reading frames; it largely depends on the length and translation level of the intORF relative to the overlapping canonical CDS.
Long noncoding RNA ORFs (lncRNA-ORFs)	Translated sequences located within transcripts currently annotated as long noncoding RNAs (lncRNAs), including long intervening/intergenic noncoding RNAs (lincRNAs), long noncoding RNAs that host small RNA species (encompassing microRNAs, snoRNAs, etc.), antisense RNAs and others.	May produce independently functional proteins. Typically lack strong sequence conservation.

observing that 3,085 of 7,264 Ribo-seq ORFs were found by more than one publication (Supplementary Fig. 2; Supplementary Tables 2 and 3). However, although such 'reproducibility' can demonstrate consistency in Ribo-seq signal, it neither provides insights into biological function nor indicates that the 4,179 non-replicated ORFs are 'false'. A major goal of Phase II will be to incorporate a greater diversity of

human cell types and tissues for improved estimates of ORF reproducibility, expression patterns and potential cell type specificity, along with further evaluation of criteria to quantify the technical confidence in Ribo-seq ORF calls.

Furthermore, Phase I excluded many translations by restricting the consensus set to ATG-initiated 'cognate' translations of at least 16 aa in length. Although these

tiny ORFs may provoke skepticism in the absence of additional evidence—the smallest annotated human protein is 24 aa—there may be no lower size limit for a functional ORF¹⁵. For example, the tarsal-less (*tal*) gene produces a polycistronic transcript translated into proteins as short as 11 aa in several insect species¹⁶. Furthermore, the inclusion of ORFs initiated with near-cognate start codons can be complicated by ambiguous predictions of initiation site positions¹⁷. Ribo-seq following treatment with lactimidomycin or homoharringtonine, which inhibit translation elongation and result in accumulation of sequencing reads at the putative start sites, can help to identify near-cognate start sites^{17,18}. Such datasets will be leveraged by our future Phase II efforts. For our current annotation resource, we have separately aggregated the Ribo-seq ORFs with near-cognate start codons or translations shorter than 16 codons (Supplementary Fig. 3a–c and Supplementary Tables 4 and 5), rather than including them in the Phase I catalog.

A core aim of Phase II will be to identify which Ribo-seq ORFs participate in cell physiology and how they do so. One aspect is distinguishing between cellular function mediated by a stable protein and functionality imparted at the level of translation itself. We here use 'protein' as an umbrella term for protein, peptide and polypeptide, although we recognize that the terms polypeptide, micropeptide or microprotein are commonly used for small protein molecules (Table 2). Because of the challenges of protein sequencing, evolutionary analysis has played a major historical role in ORF annotation, which is based on the assumption that the evolution of translated sequences is driven by selection at the protein level. Within our Phase I dataset, 75 Phase I replicated Ribo-seq ORFs (2.4%) present evidence of potential protein-level constraint as measured by PhyloCSF¹⁹ (Supplementary Fig. 3d–f); among these, ten have now been classified as protein coding by GENCODE (Supplementary Table 6).

Nonetheless, the evolutionary profile of many Phase I Ribo-seq ORFs remains hard to interpret. In part, this is because distinguishing ORF selection at the protein and DNA levels can be especially difficult for very small regions, and Ribo-seq ORFs are typically much smaller than those of known annotated proteins (Supplementary Fig. 3g–j). A second drawback is that evolutionary analysis cannot infer the protein-coding or regulatory potential of evolutionarily 'young' de novo Ribo-seq ORFs²⁰. Reference annotation projects remain skeptical

about the existence of proteins that are not deeply conserved, despite the fact that some young proteins clearly do participate in cellular physiology^{20,21}. Furthermore, there is a substantial knowledge gap in regard to the mode and tempo of regulatory ORF evolution. Here, genetic variation within human populations may provide insights. For example, Whiffin et al.²² recently used the gnomAD human variation dataset to identify 3,191 genes in which uORF-perturbing variants are likely to be deleterious, thereby inferring the physiological importance of these translations. Meanwhile, Neville et al.²³ used the same dataset to find aggregate evidence of selective pressure against deleterious variants in their nORFs.org catalog¹³, which is especially pronounced for STOP-gain variants in uORFs. In prostate cancer, a recent analysis of 5' UTR variants found regulatory roles for several uORFs²³.

Although Ribo-seq ORFs may have regulatory roles irrespective of an encoded protein, the first step in confirming a protein-level physiological role for such an ORF is to demonstrate the existence of the protein in the cell. MS is a widely accepted approach to catalog the proteome, and its utility will be an important area of investigation for Phase II. At present, 609 of 7,264 Ribo-seq ORFs have been reported to have support in published MS datasets (Supplementary Table 10). However, different groups use distinct methodologies and parameters for MS, and for Phase I these findings are simply reported in Supplementary Tables 2 and 3 without further investigation. Reference annotation projects have historically favored high-stringency MS approaches, and the Human Proteome Organization (HUPO)/ Human Proteome Project (HPP)—which aims to produce a full annotation of the human proteome—has published guidelines to standardize the nature of MS evidence required to annotate a human protein²⁴. As one facet of our development of an MS workflow, these Ribo-seq ORFs have been added to the PeptideAtlas analytical pipeline, which is used by HUPO. In Phase II, our projects will jointly examine the question of how best to use MS data to define which Ribo-seq ORFs produce proteins. For reference annotation, we see two aspects to this: first, how to set standards for accepting and reporting potential MS support for a prospective Ribo-seq ORF protein; and second, how to define the point at which the body of evidence supports protein-coding annotation.

These aspects are illustrated by a preliminary analysis, which took advantage

of the fact that 333 of our Ribo-seq ORFs are present in sequences previously queried by the PeptideAtlas workflow (Supplementary Methods). We find single-mapping peptide-spectrum matches (PSMs) for 13 Ribo-seq ORFs (Supplementary Table 11); all but one are supported by a single PSM each, whereas most of the peptides identified are not fully tryptic (two examples are presented in Supplementary Fig. 4). The majority of observed PSMs derive from human leukocyte antigen (HLA) peptidome datasets, which is consistent with prior proteomic analyses demonstrating enrichment for peptides mapping to Ribo-seq ORFs in immunopeptidome data^{25–27}. We emphasize that this preliminary analysis was not a full remapping of MS data and involved only a fraction of the Ribo-seq ORFs; a larger, focused effort will be forthcoming.

There are multiple causes contributing to the fact that Ribo-seq ORFs and certain classes of canonical proteins are infrequently detected in MS data, which are summarized elsewhere²⁸. One consideration for HUPO is that an MS-based ‘canonical’ protein assignment requires multiple PSMs, ideally based on non-overlapping tryptic peptides. Although we recognize the value of these guidelines, very small proteins may be ‘less discoverable’ by MS, especially due to a paucity of identifiable tryptic fragments²⁸. Notably, nearly 1,500 protein-coding genes annotated by GENCODE, UniProt and HGNC do not presently have MS support recognized by HUPO²⁴. Moving forward, we are committed to examining all potential protein-coding Ribo-seq ORF cases with full manual gene annotation processes, and we plan to expand this workflow to include manual analysis of the peptide spectra by PeptideAtlas.

Although the value of MS in identifying translated proteins is indisputable, we believe a broader ‘gold standard’ for evidence should employ additional methodologies, such as epitope tagging combined with western blot imaging or endogenous antibody work; HUPO already incorporates such data in collaboration with the Human Protein Atlas²⁴. Consideration also needs to be given to emerging proteomics technologies, such as targeted proteomics workflows and immunopeptidomics, and progress is being made in medium-throughput functional screening assays. For example, recent large-scale studies have translated hundreds of Ribo-seq ORFs in mammalian cells through exogenous expression, finding that nearly 50% may stably produce proteins, despite little evidence of evolutionary constraint^{2,6,27}.

In addition to their evaluation as proteins or regulatory units, the reference annotation of Ribo-seq ORFs necessitates the creation of integrated workflows to interpret overlapping variants, and notwithstanding great community interest in this field, standardized approaches are not yet available. We emphasize that variant interpretation pipelines designed to classify CDS mutations may be unsuitable for Ribo-seq ORFs (Table 1), and that a minority of overlapping variants fall within sequences displaying amino-acid-level constraint. Neville et al.¹³ found that their nORFs.org catalog contains 48 Human Gene Mutation Database or ClinVar variants that are already considered pathogenic or likely to be pathogenic, even though they do not disrupt annotated CDSs. Although these variants may affect noncanonical ORFs, it will be important to define their mechanisms of action through experimental studies, as alternative explanations for pathogenicity, such as the creation of cryptic splice sites, are supported in certain cases. After exclusion of variants in Ribo-seq ORFs that overlap annotated CDSs, a total of 1,142 single-nucleotide variants present in the ClinVar database²⁹ were located within our aggregated set of Phase I Ribo-Seq ORFs (Supplementary Methods). Fewer than 2% of these variants have been classified as pathogenic or likely to be pathogenic, but this is likely to be an underestimate because the absence of pathogenesis is commonly inferred from the absence of overlap with known coding features, and because ClinVar variant coverage is heavily skewed toward annotated CDSs.

Furthermore, there is major interest in the application of Ribo-seq to study human disease. In particular, it is being widely used to explore the dynamics of translation in cancer cells with aberrant proteins as diagnostic markers or targets for immunotherapy^{25,26,30}. At present, reference annotation projects do not attempt to distinguish aberrant translation events from those that contribute to ‘normal’ physiology. It will be important to deduce the fraction of Ribo-seq ORFs that encode proteins that exist in normal cellular conditions. Conversely, we envisage the value of classifying potentially aberrant translations within Phase II through a distinct annotation framework.

Our intention is for the Ribo-seq Phase I catalog to be seen as a pragmatic interim solution to a long-term problem. We believe that reference annotation databases can advance both scientific and clinical research through the propagation and standardization of Ribo-seq ORF datasets, even—and perhaps especially—while

the phenotypic impact of these features remains uncertain. As biological knowledge improves, this will support the development of more accurate annotations and variant interpretations, with the potential to yield substantial insights across all aspects of human biology. In this spirit, we hope the results of Phase I of this project will be useful and beneficial to the community and invite interested labs to join our future Phase II efforts.

Endorsement

HUPO/HPP Executive Committee members affirming support for this work are Rudolf Aebbersold, Cecilia Lindskog Bergström, Yu-Ju Chen, Fernando Corrales, Lydie Lane, Siqi Liu, Edward Nice, Gilbert Omenn, Christopher Overall, Young-Ki Paik, Charles Pineau, Michael Roehrl and Susan Weintraub. This work is further endorsed by Piero Carninci from Human Technopole and RIKEN.

Jonathan M. Mudge^{1,4,22}, Jorge Ruiz-Orera^{2,4,2}, John R. Prensner^{3,4,5,42}, Marie A. Brunet⁶, Ferriol Calvet¹, Irwin Jungreis^{3,7}, Jose Manuel Gonzalez¹, Michele Magrane¹, Thomas F. Martinez^{8,9}, Jana Felicitas Schulz¹⁰, Yucheng T. Yang^{10,11}, M. Mar Albà^{12,13}, Julie L. Aspden^{14,15}, Pavel V. Baranov¹⁶, Ariel A. Bazzini^{17,18}, Elspeth Bruford^{19,19}, Maria Jesus Martin¹, Lorenzo Calviello^{20,21}, Anne-Ruxandra Carvunis^{22,23}, Jin Chen²⁴, Juan Pablo Couso²⁵, Eric W. Deutsch²⁶, Paul Flicek²⁷, Adam Frankish¹, Mark Gerstein^{10,27,28,29}, Norbert Hubner^{10,30,31}, Nicholas T. Ingolia³², Manolis Kellis^{3,7}, Gerben Menschaert³³, Robert L. Moritz²⁶, Uwe Ohler^{34,35,36}, Xavier Roucou³⁷, Alan Saghatelian³⁸, Jonathan S. Weissman^{38,39,40} and Sebastiaan van Heesch^{41,42}

¹European Molecular Biology Laboratory, European Bioinformatics Institute, Wellcome Genome Campus, Hinxton, Cambridge, UK. ²Cardiovascular and Metabolic Sciences, Max Delbrück Center for Molecular Medicine in the Helmholtz Association (MDC), Berlin, Germany. ³Broad Institute of MIT and Harvard, Cambridge, MA, USA. ⁴Department of Pediatric Oncology, Dana-Farber Cancer Institute, Boston, MA, USA. ⁵Division of Pediatric Hematology/Oncology, Boston Children's Hospital, Boston, MA, USA. ⁶Department of Pediatrics, Medical Genetics Service, Université de Sherbrooke, Sherbrooke, Quebec, Canada. ⁷MIT Computer Science and Artificial Intelligence Laboratory, Cambridge, MA, USA. ⁸Clayton Foundation Laboratories for Peptide Biology, Salk Institute for Biological Studies, La Jolla, CA, USA. ⁹Department

of Pharmaceutical Sciences, University of California, Irvine, CA, USA. ¹⁰Program in Computational Biology & Bioinformatics, Yale University, New Haven, CT, USA. ¹¹Department of Molecular Biophysics & Biochemistry, Yale University, New Haven, CT, USA. ¹²Evolutionary Genomics Group, Research Programme on Biomedical Informatics, Hospital del Mar Research Institute (IMIM) and Universitat Pompeu Fabra (UPF), Barcelona, Spain. ¹³Catalan Institution for Research and Advanced Studies (ICREA), Barcelona, Spain.

¹⁴School of Molecular and Cellular Biology, Faculty of Biological Sciences, University of Leeds, Leeds, UK. ¹⁵LeedsOmics, University of Leeds, Leeds, UK. ¹⁶School of Biochemistry and Cell Biology, University College Cork, Cork, Ireland. ¹⁷Stowers Institute for Medical Research, Kansas City, MO, USA. ¹⁸Department of Molecular and Integrative Physiology, University of Kansas Medical Center, Kansas City, KS, USA. ¹⁹Department of Haematology, University of Cambridge School of Clinical Medicine, Cambridge, UK. ²⁰Functional Genomics Centre, Human Technopole, Milan, Italy. ²¹Computational Biology Centre, Human Technopole, Milan, Italy. ²²Department of Computational and Systems Biology, School of Medicine, University of Pittsburgh, Pittsburgh, PA, USA. ²³Pittsburgh Center for Evolutionary Biology and Medicine, School of Medicine, University of Pittsburgh, Pittsburgh, PA, USA. ²⁴Department of Pharmacology and Cecil H. and Ida Green Center for Reproductive Biology Sciences, University of Texas Southwestern Medical Center, Dallas, TX, USA. ²⁵Centro Andaluz de Biología del Desarrollo, CSIC-UPO, Seville, Spain. ²⁶Institute for Systems Biology, Seattle, WA, USA. ²⁷Department of Molecular Biophysics and Biochemistry, Yale University, New Haven, CT, USA.

²⁸Department of Computer Science, Yale University, New Haven, CT, USA. ²⁹Department of Statistics & Data Science, Yale University, New Haven, CT, USA. ³⁰Charité-Universitätsmedizin, Berlin, Germany. ³¹DZHK (German Centre for Cardiovascular Research), Partner Site Berlin, Berlin, Germany. ³²Department of Molecular and Cell Biology and California Institute for Quantitative Biosciences, University of California, Berkeley, Berkeley, CA, USA. ³³Biobix, Lab of Bioinformatics and Computational Genomics, Department of Mathematical Modelling, Statistics and Bioinformatics, Ghent University, Ghent, Belgium. ³⁴Berlin Institute for Medical Systems Biology, Max Delbrück Center for Molecular Medicine in the Helmholtz Association (MDC), Berlin, Germany. ³⁵Department of Biology, Humboldt-Universität zu Berlin, Berlin, Germany. ³⁶Department of Computer Science, Humboldt-Universität zu Berlin, Berlin, Germany. ³⁷Department of Biochemistry and Functional Genomics, Université de Sherbrooke, Sherbrooke, Québec, Canada. ³⁸Department of Biology, Massachusetts Institute of Technology, Cambridge, MA, USA. ³⁹Whitehead Institute for Biomedical Research, Cambridge, MA, USA. ⁴⁰Howard Hughes Medical Institute, Massachusetts Institute of Technology, Cambridge, MA, USA. ⁴¹Princess Máxima Center for Pediatric

Oncology, Utrecht, the Netherlands. ⁴²These authors contributed equally: Jonathan M. Mudge, Jorge Ruiz-Orera, John R. Prensner, Sebastiaan van Heesch ✉e-mail: jmudge@ebi.ac.uk; jorge.ruizorera@mdc-berlin.de; prensner@broadinstitute.org; S.vanHeesch@prinsesmaximacentrum.nl

Published online: 13 July 2022
<https://doi.org/10.1038/s41587-022-01369-0>

References

- Ingolia, N. T., Ghaemmaghamsi, S., Newman, J. R. S. & Weissman, J. S. *Science* **324**, 218–223 (2009).
- van Heesch, S. et al. *Cell* **178**, 242–260.e29 (2019).
- Ji, Z., Song, R., Regev, A. & Struhl, K. *eLife* **4**, e08890 (2015).
- Calviello, L. et al. *Nat. Methods* **13**, 165–170 (2016).
- Martinez, T. F. et al. *Nat. Chem. Biol.* **16**, 458–468 (2020).
- Chen, J. et al. *Science* **367**, 1140–1146 (2020).
- Gaertner, B. et al. *eLife* **9**, e58659 (2020).
- Raj, A. et al. *eLife* **5**, e13328 (2016).
- Olexiouk, V., Van Criekeing, W. & Menschaert, G. *Nucleic Acids Res.* **46**, D497–D502 (2018). D1.
- Michel, A. M., Kinary, S. J., O'Connor, P. B. F., Mullan, J. P. & Baranov, P. V. *Nucleic Acids Res.* **46**, D823–D830 (2018). D1.
- Kinary, S. J., O'Connor, P. B. F., Michel, A. M. & Baranov, P. V. *Nucleic Acids Res.* **47**, D847–D852 (2019). D1.
- Brunet, M. A. et al. *Nucleic Acids Res.* **47**, D403–D410 (2019). D1.
- Neville, M. D. C. et al. *Genome Res.* **31**, 327–336 (2020).
- McGillivray, P. et al. *Nucleic Acids Res.* **46**, 3326–3338 (2018).
- Vattem, K. M. & Wek, R. C. *Proc. Natl Acad. Sci. USA* **101**, 11269–11274 (2004).
- Galindo, M. I., Pueyo, J. I., Fouix, S., Bishop, S. A. & Couso, J. P. *PLoS Biol.* **5**, e106 (2007).
- Lee, S. et al. *Proc. Natl Acad. Sci. USA* **109**, E2424–E2432 (2012).
- Ingolia, N. T., Lareau, L. F. & Weissman, J. S. *Cell* **147**, 789–802 (2011).
- Lin, M. F., Jungreis, I. & Kellis, M. *Bioinformatics* **27**, i275–i282 (2011).
- Levy, A. *Nature* **574**, 314–316 (2019).
- Ruiz-Orera, J. & Albà, M. M. *Trends Genet.* **35**, 186–198 (2019).
- Whiffin, N. et al. *Nat. Commun.* **11**, 1 (2020).
- Lim, Y. et al. *Nat. Commun.* **12**, 4217 (2021).
- Adhikari, S. et al. *Nat. Commun.* **11**, 5301 (2020).
- Ouspenskaia, T. et al. *Nat. Biotechnol.* **40**, 209–217 (2022).
- Laumont, C. M. et al. *Sci. Transl. Med.* **10**, eaau5516 (2018).
- Prensner, J. R. et al. *Nat. Biotechnol.* **39**, 697–704 (2021).
- Omenn, G. S., Lane, L., Lundberg, E. K., Overall, C. M. & Deutsch, E. W. *J. Proteome Res.* **16**, 4281–4287 (2017).
- Landrum, M. J. et al. *Nucleic Acids Res.* **48**, D835–D844 (2020). D1.
- Chong, C. et al. *Nat. Commun.* **11**, 1293 (2020).

Acknowledgements

A.F., J.M.M., F.C. and P.F. are supported by the Wellcome Trust (grant number 108749/Z/15/Z), the National Human Genome Research Institute (NHGRI) of the US National Institutes of Health (NIH) under award number 2U41HG007234 and the European Molecular Biology Laboratory (EMBL). For the purpose of open access, the author has applied a CC BY public copyright license to any Author Accepted Manuscript version arising from this submission. The content is solely the responsibility of the authors and does not necessarily represent the official views of the National Institutes of Health. Ensembl is a registered trademark of EMBL. M.G. and Y.T.Y. are supported by NHGRI, NIH, under award number 2U41HG007234. I.J. and M.K. are supported by National Human Genome Research Institute of the NIH under award numbers 2U41HG007234 and R01 HG004037. UniProt is supported by the NHGRI, NIH, under award number (U24HG007822), EMBL core funds and the Swiss Federal Government through the State Secretariat for Education, Research and Innovation (SERI). J.R.P. is supported by the Harvard K-12 in Central Nervous System tumors (5K12 CA 90354-18), the Alex's Lemonade Stand Foundation Young Investigator Award (no. 21-23983) and the Musella Foundation for Brain Tumor Research. T.F.M. is supported by the NIH under award number F32GM123685. M.M.A. acknowledges

funding from the Spanish Government grant PGC2018-094091-B-I00 (MCI/AEI/FEDER, EU) and AGAUR grant 2017SGR01020. J.C. is supported by the NIH Pathway to Independence Award (R00 GM134154) and the Cancer Prevent and Research Institute of Texas (RR200095). J.S.W. is supported by the Howard Hughes Medical Institute. A.A.B. is supported by the Stowers Institute for Medical Research and the NIH (R01 GM136849). A.-R.C. is supported by funds provided by the Searle Scholars program, the Sloan Research Fellowship in Computational and Evolutionary Molecular Biology and the National Institute of General Medical Sciences, NIH, award number DP2GM137422. P.V.B. wishes to acknowledge the support from the Investigator in Science Award (grant number 210692/Z/18/Z) by SFI-HRB-Wellcome Trust Biomedical Research Partnership and from Russian Science Foundation (grant number 20-14-00121). N.H. is the recipient of a European Research Council advanced grant under the European Union's Horizon 2020 research and innovation programme (grant agreement no. AdG788970). N.H. is supported by a grant from the Leducq Foundation (11 CVD-01). The work of the HGNC is funded by the Wellcome Trust (208349/Z/17/Z) and the NHGRI, NIH (under award number U24HG003345). X.R.

and M.A.B. are funded by the Canadian Institutes of Health Research, grant PJT-175322. X.R. is funded as a Canada Research Chair in functional proteomics and discovery of novel proteins. N.T.I. is supported by the NIH under award number R01 GM130996. R.L.M. and E.W.D. are supported by NIH grants R01GM087221, R24GM127667, U19AG023122, 1S10OD026936-01 and US National Science Foundation grant DBI-1933311. J.L.A. is supported by the Medical Research Council (MR/N000471/1). M.A.B. is supported by a Junior 1 fellowship from the Fonds de Recherche du Québec-Santé.

Author contributions

J.M.M., J.R.-O., J.R.P. and S.v.H. conceptualized the work and supervised the international collaboration. J.R.-O., J.M.G., M.M., M.J.M., F.C., E.B., E.W.D., R.L.M. and J.M.M. performed data curation. All authors contributed to standardization of the data analysis approach. All authors contributed to discussions on Phase I and II of this effort and continue to provide scientific oversight. A.F., P.F., M.J.M., G.M., Y.T.Y., J.R.P., T.F.M., M.M.A., J.C., J.S.W., A.A.B., A.-R.C., P.V.B., N.H., X.R., M.A.B., N.T.I., E.B., E.W.D. and R.L.M. provided funding. J.M.M., J.R.-O., J.R.P.

and S.v.H. wrote the original manuscript draft. All authors reviewed the manuscript and provided edits. All authors approved the final manuscript.

Competing interests

P.V.B. is a co-founder of RiboMaps Ltd., which provides Ribo-seq analysis as a commercial service, including identification of translated ORFs. A.R.C. is a member of the scientific advisory board for Flagship Labs 69, Inc. P.F. is a member of the scientific advisory boards of Fabric Genomics, Inc., and Eagle Genomics, Ltd. The other authors declare no competing interests.

Additional information

Supplementary information The online version contains supplementary material available at <https://doi.org/10.1038/s41587-022-01369-0>.

Peer review information *Nature Biotechnology* thanks Sudhakaran Prabakaran, Mathias Uhlén and the other, anonymous, reviewer(s) for their contribution to the peer review of this work.



PRO-ACTive sharing of clinical data

To the Editor — The open sharing of clinical data for research poses challenges not only in resolving consent, privacy and intellectual property issues associated with trial results^{1,2}, but also in subsequently facilitating access to and utilization of the data. Controlled-access systems can place onerous restrictions on industry-based researchers, require arduous application processes and involve long review or authorization times for users. Indeed, analyses of major efforts, such as ClinicalStudyDataRequest.com (CSDR)^{3,4} and the UK Health Research Authority's (HRA) Assessment Review Portal (HARP) database⁵, suggest that fewer than 50% of the available trial datasets in these resources have been accessed and analyzed by researchers after launch. Here, we outline several strategies that ensured that researcher-engagement with an open-access ALS clinical trial data resource reached its full potential. We hope that our insights will be instructive for others seeking to galvanize open clinical data sharing efforts within the broader research community.

Prize4Life, a non-profit organization focused on accelerating treatments and a cure for ALS, created the Pooled Resource Open-Access ALS Clinical Trials (PRO-ACT) database (<https://ncrl.partners.org/ProACT/Document/DisplayLatest/9>) in collaboration with Massachusetts General Hospital's Neurological Clinical Research Institute and with funding from the ALS Therapy Alliance. PRO-ACT was publicly launched in 2012 (with further data incorporated in 2015), including

Table 1 | Comparative access and usage for PRO-ACT and other open databases

Aspect	PRO-ACT	CSDR and HARP ³⁻⁵
Review time	-1 day	6-12 months
Number of requests	>2,500	300-500
Percentage of requests accepted	87%	<50%
Access requirements	<ul style="list-style-type: none"> Brief terms and conditions (for example, privacy and data integrity). Short summary of research plan 	<ul style="list-style-type: none"> Institutional ethics approval Comprehensive research plan Formal committee review
Cost of access	Free to access the data; inexpensive to provide the data	Time- and resource-intensive to access the data
Data included	<ul style="list-style-type: none"> Placebo and active Entire dataset downloadable 	<ul style="list-style-type: none"> Placebo data Frequent requirement to only use data within a cloud-based or managed environment
Efforts to publicize	<ul style="list-style-type: none"> Use of the prize model to attract widespread interest across disciplines Use of patients to explain and personify research needs 	Typically limited, and discipline based only
Number of publications	>80	<10 across examples ⁴⁻⁶

demographics, family history, medical history (including use of frontline treatment riluzole), vital capacity, adverse events and other data types from both the placebo and active arms of over 20 clinical trials^{6,7}. The database currently holds >10,000 fully deidentified ALS patient records from 23 phase 2 and 3 clinical trials, representing the largest aggregation of publicly available ALS clinical data.

In contrast to other clinical trial repositories (Table 1), PRO-ACT has been

widely accessed by >2,500 users, from >50 countries, including dozens of universities, several governmental agencies and >50 drug development companies. Over 80 publications have used PRO-ACT as a primary data source. Last year, the success of PRO-ACT and its value to the ALS field was acknowledged when the database's creators received the Healey Center Prize for Innovation in ALS⁸. Several factors have contributed to the success of PRO-ACT in engaging the wider research community.