

Perspectives on ENCODE

<https://doi.org/10.1038/s41586-020-2449-8>

Received: 7 December 2019

Accepted: 5 May 2020

Published online: 29 July 2020

Open access

 Check for updates

The ENCODE Project Consortium*, Michael P. Snyder^{1,2✉}, Thomas R. Gingeras³, Jill E. Moore⁴, Zhiping Weng^{4,5,6}, Mark B. Gerstein⁷, Bing Ren^{8,9}, Ross C. Hardison¹⁰, John A. Stamatoyannopoulos^{11,12,13}, Brenton R. Graveley¹⁴, Elise A. Feingold¹⁵, Michael J. Pazin¹⁵, Michael Pagan¹⁵, Daniel A. Gilchrist¹⁵, Benjamin C. Hitz¹, J. Michael Cherry¹, Bradley E. Bernstein¹⁶, Eric M. Mendenhall^{17,18}, Daniel R. Zerbino¹⁹, Adam Frankish¹⁹, Paul Flicek¹⁹ & Richard M. Myers¹⁸

The Encyclopedia of DNA Elements (ENCODE) Project launched in 2003 with the long-term goal of developing a comprehensive map of functional elements in the human genome. These included genes, biochemical regions associated with gene regulation (for example, transcription factor binding sites, open chromatin, and histone marks) and transcript isoforms. The marks serve as sites for candidate *cis*-regulatory elements (cCREs) that may serve functional roles in regulating gene expression¹. The project has been extended to model organisms, particularly the mouse. In the third phase of ENCODE, nearly a million and more than 300,000 cCRE annotations have been generated for human and mouse, respectively, and these have provided a valuable resource for the scientific community.

The ENCODE Project was launched in 2003, as the first nearly complete human genome sequence was reported². At that time, our understanding of the human genome was limited. For example, although 5% of the genome was known to be under purifying selection in placental mammals^{3,4}, our knowledge of specific elements, particularly with regards to non-protein coding genes and regulatory regions, was restricted to a few well-studied loci^{2,5}.

ENCODE commenced as an ambitious effort to comprehensively annotate the elements in the human genome, such as genes, control elements, and transcript isoforms, and was later expanded to annotate the genomes of several model organisms. Mapping assays identified biochemical activities and thus candidate regulatory elements.

Analyses of the human genome in ENCODE proceeded in successive phases (Extended Data Fig. 1). Phase I (2003–2007) interrogated a specified 1% of the human genome in order to evaluate emerging technologies⁶. Half of this 1% was in regions of high interest, and the other half was chosen to sample the range of genomic features (such as G+C content and genes). Microarray-based assays were used to map transcribed regions, open chromatin, and regions associated with transcription factors and histone modification in a wide variety of cell lines, and these assays began to reveal the basic organizational features of the human genome and transcriptome. Phase II (2007–2012) introduced sequencing-based technologies (for example, chromatin immunoprecipitation with sequencing (ChIP-seq) and RNA sequencing (RNA-seq)) that interrogated the whole human genome and transcriptome⁷. General assays such as transcript, open-chromatin and histone

modification mapping were used on a wide variety of cell lines, while more specific assays, such as mapping transcription factor binding regions, were performed extensively on a smaller number of cell lines to provide detailed annotations on, and to investigate the relationships of, many regulatory proteins across the genome. Transcriptome analysis of subcellular compartments (the nucleus, cytosol and subnuclear compartments) of these cells enabled the locations of transcripts to be analysed⁷.

ENCODE phase III

ENCODE 3 (2012–2017) expanded production and added new types of assays⁸ (Fig. 1, Extended Data Fig. 1), which revealed landscapes of RNA binding and the 3D organization of chromatin via methods such as chromatin interaction analysis by paired-end tagging (ChIA-PET) and Hi-C chromosome conformation capture. Phases 2 and 3 delivered 9,239 experiments (7,495 in human and 1,744 in mouse) in more than 500 cell types and tissues, including mapping of transcribed regions and transcript isoforms, regions of transcripts recognized by RNA-binding proteins, transcription factor binding regions, and regions that harbour specific histone modifications, open chromatin, and 3D chromatin interactions. The results of all of these experiments are available at the ENCODE portal (<http://www.encodeproject.org>). These efforts, combined with those of related projects and many other laboratories, have produced a greatly enhanced view of the human genome (Fig. 2), identifying 20,225 protein-coding and 37,595 noncoding genes

¹Department of Genetics, School of Medicine, Stanford University, Palo Alto, CA, USA. ²Cardiovascular Institute, Stanford School of Medicine, Stanford, CA, USA. ³Functional Genomics, Cold Spring Harbor Laboratory, Cold Spring Harbor, NY, USA. ⁴University of Massachusetts Medical School, Program in Bioinformatics and Integrative Biology, Worcester, MA, USA. ⁵Department of Thoracic Surgery, Clinical Translational Research Center, Shanghai Pulmonary Hospital, The School of Life Sciences and Technology, Tongji University, Shanghai, China. ⁶Bioinformatics Program, Boston University, Boston, MA, USA. ⁷Yale University, New Haven, CT, USA. ⁸Ludwig Institute for Cancer Research, University of California, San Diego, La Jolla, CA, USA. ⁹Center for Epigenomics, University of California, San Diego, La Jolla, CA, USA. ¹⁰Department of Biochemistry and Molecular Biology, The Pennsylvania State University, University Park, PA, USA. ¹¹Altius Institute for Biomedical Sciences, Seattle, WA, USA. ¹²Department of Genome Sciences, University of Washington, Seattle, WA, USA. ¹³Department of Medicine, University of Washington, Seattle, WA, USA. ¹⁴Department of Genetics and Genome Sciences, Institute for Systems Genomics, UConn Health, Farmington, CT, USA. ¹⁵National Human Genome Research Institute, National Institutes of Health, Bethesda, MD, USA. ¹⁶Broad Institute and Department of Pathology, Massachusetts General Hospital and Harvard Medical School, Boston, MA, USA. ¹⁷Biological Sciences, University of Alabama in Huntsville, Huntsville, AL, USA. ¹⁸HudsonAlpha Institute for Biotechnology, Huntsville, AL, USA. ¹⁹European Molecular Biology Laboratory, European Bioinformatics Institute, Wellcome Genome Campus, Cambridge, UK. *An alphabetical list of authors and their affiliations appears online. A formatted list of authors appears in the Supplementary Information. ✉e-mail: mpsnyder@stanford.edu

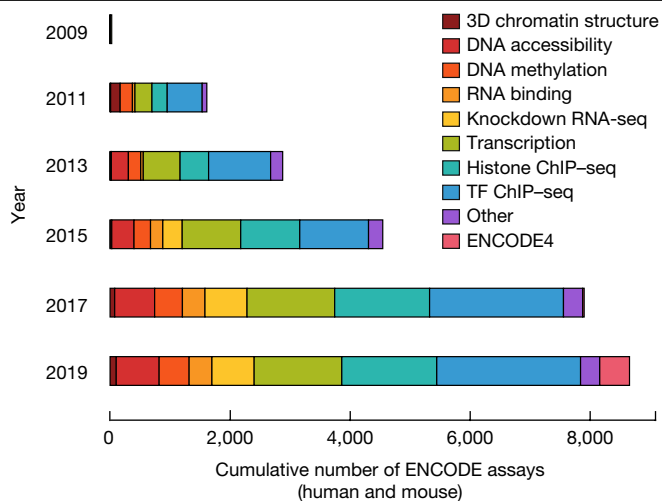


Fig. 1 | ENCODE assays by year. Accumulations of assays over the three phases of ENCODE. 3D chromatin structure includes ChIA-PET (62 experiments), Hi-C (31), and chromatin conformation capture carbon copy (5C, 13). Chromatin accessibility includes DNAase-seq (524), assay for transposase-accessible chromatin using sequencing (ATAC-seq, 129), transcription activator-like effector nuclease (TALEN)-modified DNAase-seq (40), formaldehyde-assisted isolation of regulator elements with sequencing (FAIRE-seq, 37) and micrococcal nuclease digestion with deep sequencing (MNase-seq, 2). DNA methylation includes DNAm arrays (259), WGBS (124), reduced-representation bisulfite sequencing (RRBS, 103), methylation-sensitive restriction enzyme sequencing (MRE-seq, 24) and methylated DNA immunoprecipitation coupled with next-generation sequencing (MeDIP-seq, 4). Histone modification includes ChIP-seq (1,605) on histone and modified histone targets. Knockdown transcription includes RNA-seq preceded by small interfering RNA (siRNA, 54), short hairpin RNA (shRNA, 531), clustered regularly interspaced short palindromic repeats (CRISPR, 50) or CRISPR interference (CRISPRi, 77). RNA binding includes enhanced cross-linking immunoprecipitation (eCLIP, 349), RNA bind-n-seq (158), RNA immunoprecipitation sequencing (RIP-seq, 158), RNA-binding protein immunoprecipitation-microarray profiling (RIP-chip, 32), individual nucleotide-resolution CLIP (iCLIP, 6) and Switchgear (2). Transcription includes RNA annotation and mapping of promoters for the analysis of gene expression (RAMPAGE, 155), cap analysis gene expression (CAGE, 78), RNA paired-end tag (RNA-PET, 31), microRNA-seq (114), microRNA counts (114), more classical RNA-seq (900) and RNA-microarray (170), including 112 experiments at single-cell resolution. Transcription factor (TF) binding is ChIP-seq on non-histone targets (2,443). Other assays include genotyping array (123), nascent DNA replication strand sequencing (Repli-seq, 104), replication strand arrays (Repli-chip, 63), tandem mass spectrometry (MS/MS, 14), genotyping by high-throughput sequencing (genotyping HTS, 12) and DNA-PET (6) can be looked at in detail at <https://www.encodeproject.org>.

(Fig. 2a), 2,157,387 open chromatin regions, 750,392 regions with modified histones (mono-, di- or tri-methylation of histone H3 at lysine 4 (H3K4me1, H3K4me2 or H3K4me3), or acetylation of histone 3 at lysine 27 (H3K27ac)), 1,224,154 regions bound by transcription factors and chromatin-associated proteins (Fig. 2c), 845,000 RNA subregions occupied by RNA-binding proteins, and more than 130,000 long-range interactions between chromatin loci. These annotations have greatly enhanced our view of the human genome from its original annotation in 2003 to a much richer and higher-resolution view (for example, Fig. 2d, e). Indeed, although the number of human protein-coding genes known has changed only modestly, the number of transcript isoforms, long noncoding RNAs (lncRNAs), and potential regulatory regions identified has increased greatly since the project began (Fig. 2a–c). An important part of ENCODE 3 is that the regulatory mapping efforts have now been integrated and synthesized into the first version of an encyclopedia, highlighting a registry of 0.9 million cCREs in human and 0.3 million

cCREs in mouse. Details can be found in the accompanying ENCODE paper⁸ and companion papers in this issue and other journals^{9–14}.

Technology, quality control and standards

Reaching the present annotation required a substantial expansion of technology development, from ENCODE groups and others, as well as the establishment of standards to ensure that the data are reproducible and of high quality. Most ENCODE 2 assays used sequence-based readouts (for example, RNA-seq^{15,16} and ChIP-seq^{17,18}) rather than the array-based methods^{19,20} used in the pilot phase, and in ENCODE 3, methods such as global mapping of 3D interactions¹³ and RNA-binding regions¹⁴ were added. Throughout the project, computational and visualization approaches were developed for mapping reads and integrating different data types (Supplementary Note 1).

A key feature of ENCODE is the application of data standards, including the use of independent replicates (separate experiments on two or more biological samples^{5,21}), except when precluded by the limited availability of materials (for example, postmortem human tissues). Of the 8,699 ENCODE 2 and ENCODE 3 experiments, 6,101 have independent replicates. Of equal importance was the use of well-characterized reagents, such as antibodies for mapping sites of transcription factor binding, chromatin modifications and protein–RNA interactions²². ENCODE developed protocols to test each antibody ‘lot’ to demonstrate their experimental suitability, captured extensive metadata, and implemented controlled vocabularies and ontologies. Standards for reagents, experimental data, and metadata are on the ENCODE website: <https://www.encodeproject.org/data-standards/>.

Many metrics, including sequencing depth, mapping characteristics, replicate concordance, library complexity, and signal-to-noise ratio, were used to monitor the quality of each data set, and quality thresholds were applied²¹. A minority of experiments that fell short of the standards (for example, insufficiently validated antibodies) are still reported, but are marked with a badge to indicate that an issue was found. This is a compromise for having some data versus none when an experiment did not meet ENCODE-defined thresholds.

An important component is uniform data processing. Data from the major ENCODE assays (ChIP-seq, DNase I hypersensitive sites sequencing (DNase-seq), RNA-seq, and whole-genome bisulfite sequencing (WGBS)) are uniformly processed and the processing pipelines are available for users to apply to their own data, by downloading the code from the GitHub (<http://github.com/ENCODE-DCC>) or by accessing the pipelines at the DNAnexus cloud provider. The standards and pipelines will continue to evolve as new technologies arise and are implemented.

The ENCODE Consortium is a good example of how large-scale group efforts can have a large impact on the scientific community, and many other national and international projects—including the NIH Roadmap Epigenomics Program, The Cancer Genome Atlas (TCGA), the International Human Epigenome Consortium (IHEC), BLUEPRINT, the Canadian Epigenetics, Environment and Health Research Consortium (CEEHRC), the Genotype and Tissue Expression Project (GTEx), PsychENCODE, Functional Annotation of Animal Genomes (FAANG), the Global Alliance for Genomics and Health (GA4GH), the 4D Nucleome Program (4DN), the Human Cell Atlas and the FANTOM consortium—have now formed (Supplementary Note 1). ENCODE has engaged with most of these consortia to share standards for data quality control, submission, and uniform processing and has helped to facilitate the use of common ontologies with some of these consortia. Data from the now-completed NIH Roadmap Epigenomics Program have been reprocessed and are available in the ENCODE database and are part of the Encyclopedia annotation. ENCODE continues to work with other consortia, individually and as part of the IHEC and GA4GH (for example, <http://epishare-project.org>) to increase data interoperability and the value of its resources.

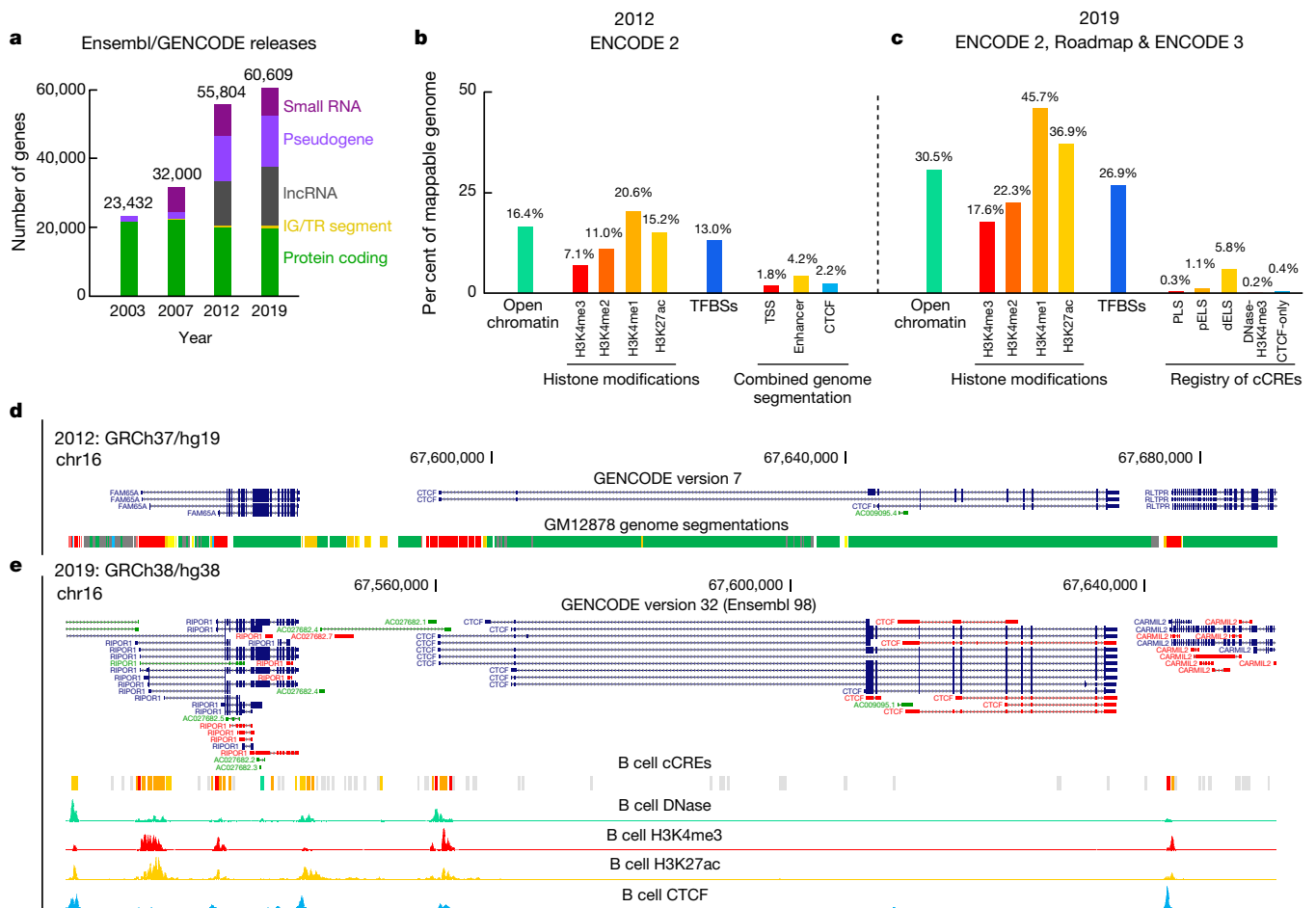


Fig. 2 | Progress in annotating the human genome. Link to high-resolution PDF file: <https://www.dropbox.com/s/rjdcrcygz15p034/perspective.pdf?dl=0>. **a**, Improvement of gene annotations in the past 15 years by GENCODE, an international gene annotation group that uses ENCODE data⁴². **b**, ENCODE annotations in 2012 with phase II data. Bars show the percentages of the mappable human genome (3.1 billion nucleotides; hg19) that were annotated as open chromatin by DNase-seq data, enriched in four types of active histone mark according to ChIP-seq data, and annotated as transcription factor binding sites (TFBSs) according to ChIP-seq data. Also shown are percentages of the genome assigned as transcription start sites (TSSs), enhancers and the insulator-binding protein (CTCF) by combining ChromHMM and Segway genome segmentations⁷. **c**, ENCODE annotations in 2019 with ENCODE 2, Roadmap, and ENCODE 3 data. The registry of cCREs developed during phase III defines 0.3%, 1.1%, 5.8%, 0.2% and 0.4% of the human

genome as cCREs with promoter-like signatures (PLS), proximal enhancer-like signatures (pELS), distal enhancer-like signatures (dELS), with high DNase, high H3K4me3 and low H3K27ac signals (DNase-H3K4me3), and bound by CTCF, respectively. **d**, A UCSC genome browser view of GENCODE genes (V7) coloured by transcript annotation (blue for coding, green for noncoding, and red for problematic) and combined genome segmentation (TSSs in red, enhancers in orange, weak enhancers in yellow, transcription in green, repressed in grey) at the *CTCF* locus on the hg19 human genome. **e**, The UCSC genome browser view of GENCODE genes (V28, coloured as in **d**) and cCREs at the *CTCF* locus on the hg38 human genome⁸. Promoter-like, enhancer-like, and CTCF-only cCREs annotated in B cells are in red, yellow, and blue, respectively. The last four tracks show the DNase, H3K4me3, H3K27ac, and CTCF signals in B cells.

ENCODE as a resource

The purpose of ENCODE is to provide valuable, accessible resources to the community. ENCODE data and derived features are available from a publicly accessible data portal (<https://www.encodeproject.org>), and consent was obtained from donors to make data freely available to the public. Raw and processed data are available directly from the cloud as an Amazon Public Data Set (<https://registry.opendata.aws/encode-project/>). The data are widely used by the scientific community—more than 2,000 publications from researchers outside of ENCODE have used ENCODE data to study diverse topics (Fig. 3). Because most disease-associated common variants are noncoding and show substantial enrichment in candidate cell-type-specific *cis* regulatory elements^{23,24}, ENCODE-derived resources, both in isolation and in conjunction with data from other resources (for example, GTEx), can help to identify and interpret disease-associated noncoding variants (Fig. 3a). Users engage with the data in many ways, ranging from

downloads of multiple data sets to detailed investigations of specific loci. Anyone navigating a major genome browser has access to thousands of biochemical, functional, and computational annotations to display at any genomic scale or to overlay on any sequence variant. Maps of epigenomic features relevant to gene regulation have been integrated to form a registry of discrete elements that are candidates for enhancers, promoters, or other regulatory elements. A specialized browser, SCREEN (<http://screen.encodeproject.org>), is an interface that can be used to identify and study these cCREs and associated ENCODE data and other annotations. This dynamic registry will be regularly updated as additional information is acquired.

Mouse ENCODE and modENCODE

Model organism studies have produced essential insights into almost every aspect of biology, including genome organization and function.

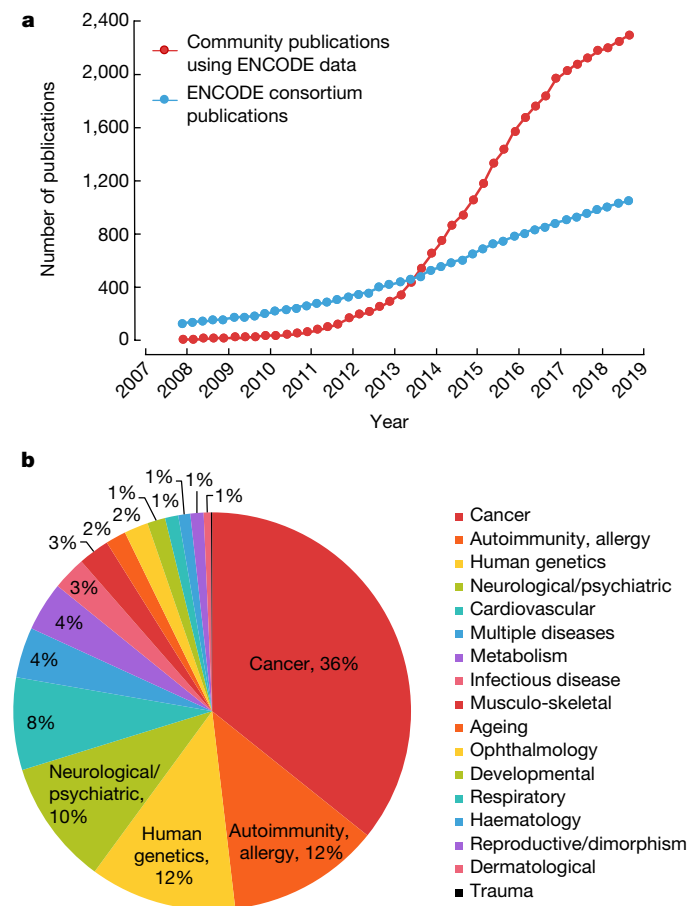


Fig. 3 | Publications using ENCODE data. The National Human Genome Research Institute (NHGRI) has identified a list of publications that used ENCODE data. This list is publicly shared to provide examples illustrating how the resource has been used (<https://www.encodeproject.org/publications/>). **a**, Publications over time. Community publications appear to use ENCODE data and do not report ENCODE grant support in PubMed; consortium publications report ENCODE grant support in PubMed. In brief, community publications are identified using two steps; first, candidates are identified through automated searches for citation of ENCODE accession numbers, ENCODE flagship papers, or resources such as HaploReg and RegulomeDB; second, candidates are manually evaluated to determine whether ENCODE data were actually used. Consortium papers are identified through automated searches of PubMed for publications that were supported at least in part by ENCODE awards, and are not further evaluated or annotated. **b**, Human disease example publications. The subset of community publications that were annotated as ‘human disease’ (other categories are basic biology, software tool, fly/worm data) were further manually categorized by disease aetiology.

During ENCODE 2, mapping of mouse epigenomic and transcriptomic features was conducted in adult mouse tissues and cell lines through the Mouse ENCODE Project²⁵, which identified 21,978 protein-coding regions, 32,168 noncoding genes, 1,192,301 open chromatin regions, 722,334 regions with modified histones H3K4me1, H3K4me2, H3K4me3, or H3K27ac, and 686,294 regions bound by transcription factors.

During ENCODE 2, a model organism ENCODE project (modENCODE^{26,27}) was conducted to characterize the transcriptome, epigenome, and transcription factor binding sites in *Drosophila melanogaster* and *Caenorhabditis elegans* tissues, developmental stages and cell lines (Extended Data Fig. 1). These organisms provided the opportunity to develop detailed records of epigenomic features and transcriptome maps throughout development, which is difficult to accomplish in humans. Deep mapping of the spatial and temporal transcriptomes

of these species has substantially enhanced the annotation of both genomes. Similarly, detailed mapping of the regulatory circuits that govern gene regulation in *Drosophila* and *C. elegans* has provided insights into general principles of genome organization and function. Mapping of transcription factor binding sites in *Drosophila* and *C. elegans* has continued after modENCODE ended in a project called model organism Encyclopedia of Regulatory Networks (modERN) and to date has characterized more than 262 transcription factors in *Drosophila* and 217 transcription factors in *C. elegans*²⁸. Collectively, the modENCODE Project has provided new insights about how the genomes of multicellular organisms direct development and maintain homeostasis.

In ENCODE phase III, experiments were carried out to characterize dynamic histone marks and accessibility, DNA methylomes, and transcriptomes in samples taken during eight mouse fetal developmental stages with up to twelve tissues per stage^{28–30} (Fig. 4). The resulting more than 1,500 datasets comprise, to our knowledge, the most comprehensive study of epigenomes and transcriptomes during the prenatal development of a mammal. Integrative analysis of these datasets has expanded our knowledge of the transcriptional regulatory networks that regulate mammalian development and underscored the role of gene regulatory mechanisms in human disease. At least 214,264 of the candidate enhancers identified in fetal mouse tissues are conserved in the human genome⁸. The human orthologues of these potential regulatory elements are significantly enriched for genetic variants that are associated with common illnesses in a tissue-restricted manner, providing information for investigations of the molecular basis of human disease^{29,30}.

The mouse data from ENCODE 3 also include the results of more than 400 experiments using transgenic reporter mice designed to assess the function of cCREs in three embryonic tissues at two developmental stages. The results of this systematic study have helped to predict the in vivo activities of cCREs. For example, stronger enrichment for epigenetic signatures of enhancer activity correlated with higher rates of validation in the corresponding tissue^{29,31}.

Finally, comparisons of epigenome and transcriptome maps across species have led to insights into the evolution of transcribed regions and regulatory information^{25,32}. Combinatorial histone modification patterns at *cis*-regulatory elements and other genomic features are broadly conserved in metazoans. These chromatin states and transcript levels are highly correlated across tissues and developmental stages in all species examined. However, a notable fraction of specific *cis*-regulatory elements undergoes sequence and functional turnover during evolution, indicating that some regulatory components show substantial plasticity in their evolution while operating in a conserved regulatory network³³.

Current limitations: phase IV and beyond

It is now apparent that elements that govern transcription, chromatin organization, splicing, and other key aspects of genome control and function are densely encoded in the human genome; however, despite the discovery of many new elements, the annotation of elements that are highly selective for particular cell types or states is lagging behind. For example, very few examples of condition-specific activation or repression of transcriptional control elements are currently annotated in ENCODE. Similarly, information from human fetal tissue, reproductive organs and primary cell types is limited. In addition, although many open chromatin regions have been mapped, the transcription factors that bind to these sequences are largely unknown, and little attention has been devoted to the analysis of repetitive sequences. Finally, although transcript heterogeneity and isoforms have been described in many cell types, full-length transcripts that represent the isoform structure of spliced exons and edits have been described for only a small number of cell types.

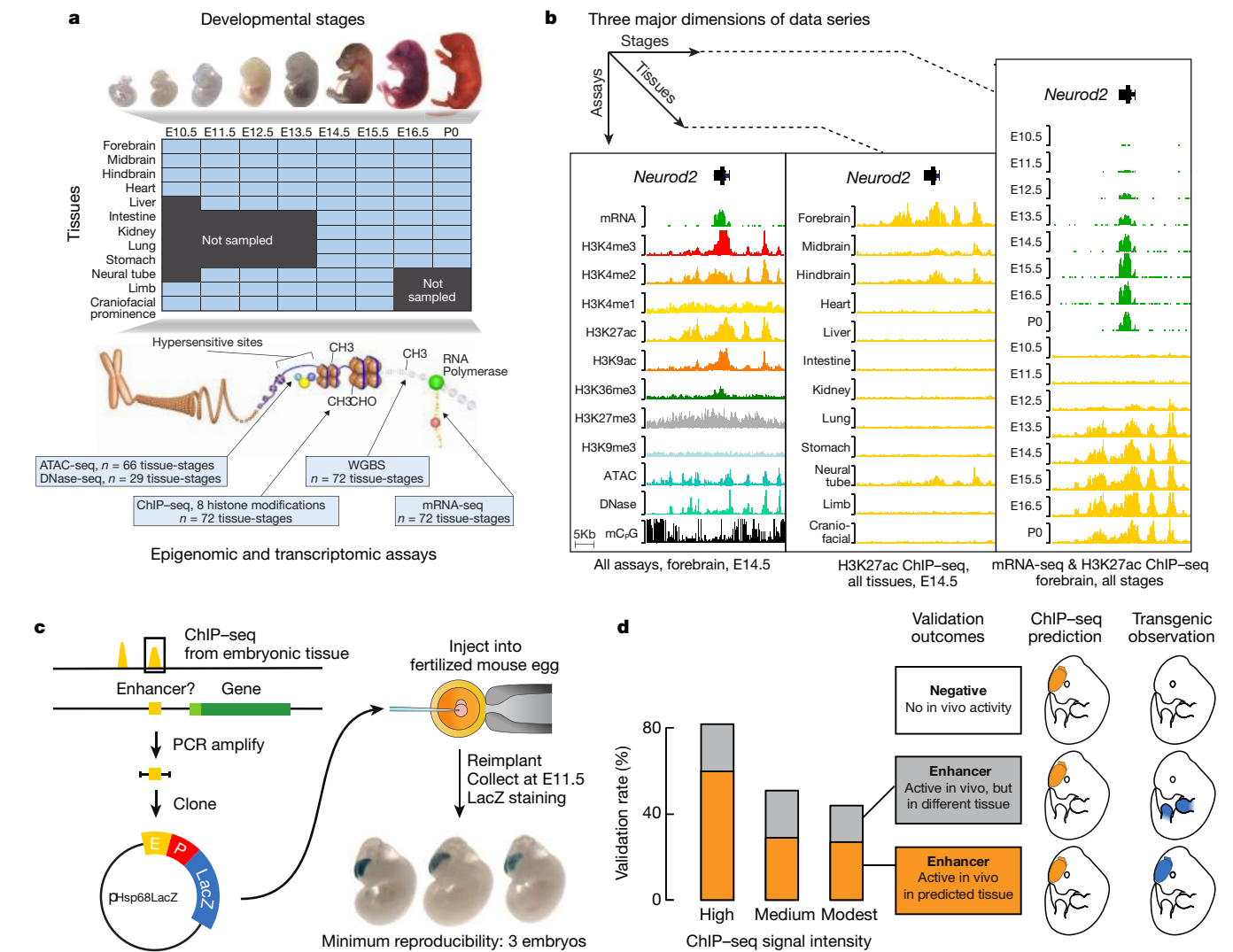


Fig. 4 | An overview of the mouse ENCODE Project in the current phase.
a, Schematic representation of ENCODE 3 mouse developmental data series. The chromatin graphic is adapted from an image by Darryl Leja (NHGRI), Ian Dunham (EBI), and M.J.P. (NHGRI). The embryo image second from the right in was adapted from ref.⁴³, an Open Access article distributed under the terms of the Creative Commons Attribution License 2.0. **b**, Three major axes of the data series: assays, tissues, and developmental stages. The region shown is chr11:98,307,637–98,344,383, mm10. **c**, A schematic diagram of the transgenic assays used to validate and characterize the function of cCREs in E11.5 and

E12.5 mouse embryos. The cCREs were selected on the basis of ChIP-seq data and cloned into a reporter vector that was then introduced into fertilized mouse eggs. The activities of the CRE were validated by tissue-specific expression patterns of the reporter gene. **d**, Results from recent transgenic assays^{8,29} to validate about 400 cCREs are summarized in a barchart, with the bars indicating the proportion of candidate CREs in each rank tier that showed reproducible reporter staining in the expected tissue (grey) or any tissue (pink).

Thus, as part of ENCODE 4, considerable effort is being devoted to expanding the cell types and tissues analysed (see URLs in Supplementary Note 1) as well as mapping the binding regions for many more transcription factors and RNA-binding proteins. These efforts are largely focused in a few reference cell lines, with the hope that improved knowledge will help with imputation or predictions in other cell states³⁴. Single-cell transcriptome capture agents³⁵ and open chromatin assays³⁶ are also being applied to increase our understanding of the cellular heterogeneity of different tissues and samples. These efforts will supplement the many related activities that are also being pursued by HCA, HuBMAP and others^{37,38}. Extensive mapping efforts of all types will continue in both the human and mouse, and parallel efforts to map transcription factor binding sites are being pursued in the *Drosophila* and *C. elegans* by the modERN Project²⁸. Full-length transcript isoforms are being elucidated in different cell types using long-read sequencing technologies³⁹. ENCODE will continue to work with other consortia, and the data from different

groups and individual laboratories will need to be consolidated into a common repository.

Importantly, although very large numbers of noncoding elements have been defined, the functional annotation of ENCODE-identified elements is still in its infancy. High-throughput reporter-based assays⁴⁰, CRISPR-based genome and epigenome editing methods⁴¹, and other high-throughput approaches are being used in the current phase of ENCODE to assess the functions of many thousands of elements and to relate those functional results to their biochemical signatures. These targeted functional assays, combined with the large-scale annotation of biochemical features, should further enhance the value of ENCODE data.

Through these and other efforts, it is expected that many more elements in the human genome will be identified across a variety of cell types and conditions, their activities will be revealed (often at the single-cell level), and their biological functions will be inferred more accurately. The development of a systems-wide understanding

of function and integration with genetic information associated with human traits will greatly enhance our understanding of human biology and disease.

Online content

Any methods, additional references, Nature Research reporting summaries, source data, extended data, supplementary information, acknowledgements, peer review information; details of author contributions and competing interests; and statements of data and code availability are available at <https://doi.org/10.1038/s41586-020-2449-8>.

- Kellis, M. et al. Defining functional DNA elements in the human genome. *Proc. Natl Acad. Sci. USA* **111**, 6131–6138 (2014).
- ENCODE Project Consortium. The ENCODE (ENCyclopedia Of DNA Elements) Project. *Science* **306**, 636–640 (2004).
- Lindblad-Toh, K. et al. A high-resolution map of human evolutionary constraint using 29 mammals. *Nature* **478**, 476–482 (2011).
- Waterston, R. H. et al. Initial sequencing and comparative analysis of the mouse genome. *Nature* **420**, 520–562 (2002).
- ENCODE Project Consortium. A user's guide to the encyclopedia of DNA elements (ENCODE). *PLoS Biol.* **9**, e1001046 (2011).
- Birney, E. et al. Identification and analysis of functional elements in 1% of the human genome by the ENCODE pilot project. *Nature* **447**, 799–816 (2007).
- ENCODE Project Consortium. An integrated encyclopedia of DNA elements in the human genome. *Nature* **489**, 57–74 (2012).
- The results of the pilot phase of ENCODE included extensive functional assays across a selected one per cent of the human genome with experiments conducted on a variety of cell lines and largely with array-based technology.
- ENCODE Project Consortium. An integrated encyclopedia of DNA elements in the human genome. *Nature* **489**, 57–74 (2012).
- The results of the second phase of ENCODE were based mostly on a large number of genome-wide assays that leveraged high-throughput sequencing technologies and were done across two 'tier one' cell lines with large-scale assays across several hundred cell and tissue types.
- The ENCODE Project Consortium et al. Expanded encyclopaedias of DNA elements in the human and mouse genomes. *Nature* <https://doi.org/10.1038/s41586-020-2493-4> (2020).
- Partridge, E. C. et al. Occupancy maps of 208 chromatin-associated proteins in one human cell type. *Nature* <https://doi.org/10.1038/s41586-020-2023-4> (2020).
- Meuleman, W. Index and biological spectrum of human DNase I hypersensitive sites. *Nature* <https://doi.org/10.1038/s41586-020-2559-3> (2020).
- Vierstra, J. et al. Global reference mapping of human transcription factor footprints. *Nature* <https://doi.org/10.1038/s41586-020-2528-x> (2020).
- Breschi, A. et al. A limited set of transcriptional programs define major cell types. Preprint at <https://doi.org/10.1101/857169> (2020).
- Grubert, F. et al. Landscape of cohesin-mediated chromatin loops in the human genome. *Nature* <https://doi.org/10.1038/s41586-020-2151-x> (2020).
- Van Nostrand, E. L. et al. A large-scale binding and functional map of human RNA binding proteins. *Nature* <https://doi.org/10.1038/s41586-020-2077-3> (2020).
- Wang, Z., Gerstein, M. & Snyder, M. RNA-seq: a revolutionary tool for transcriptomics. *Nat. Rev. Genet.* **10**, 57–63 (2009).
- Mortazavi, A., Williams, B. A., McCue, K., Schaeffer, L. & Wold, B. Mapping and quantifying mammalian transcriptomes by RNA-seq. *Nat. Methods* **5**, 621–628 (2008).
- Johnson, D. S., Mortazavi, A., Myers, R. M. & Wold, B. Genome-wide mapping of in vivo protein–DNA interactions. *Science* **316**, 1497–1502 (2007).
- Robertson, G. et al. Genome-wide profiles of STAT1 DNA association using chromatin immunoprecipitation and massively parallel sequencing. *Nat. Methods* **4**, 651–657 (2007).
- Iyer, V. R. et al. Genomic binding sites of the yeast cell-cycle transcription factors SBF and MBF. *Nature* **409**, 533–538 (2001).
- Ren, B. et al. Genome-wide location and function of DNA binding proteins. *Science* **290**, 2306–2309 (2000).
- Landt, S. G. et al. ChIP-seq guidelines and practices of the ENCODE and modENCODE consortia. *Genome Res.* **22**, 1813–1831 (2012).
- A consortium-wide effort to standardize performance, quality control and outputs of ChIP-seq experiments, including validation of antibodies, to facilitate experimental reproducibility and data utility.
- Sundararaman, B. et al. Resources for the comprehensive discovery of functional RNA elements. *Mol. Cell* **61**, 903–913 (2016).
- Maurano, M. T. et al. Systematic localization of common disease-associated variation in regulatory DNA. *Science* **337**, 1190–1195 (2012).
- Schaub, M. A., Boyle, A. P., Kundaje, A., Batzoglou, S. & Snyder, M. Linking disease associations with regulatory information in the human genome. *Genome Res.* **22**, 1748–1759 (2012).
- Yue, F. et al. A comparative encyclopedia of DNA elements in the mouse genome. *Nature* **515**, 355–364 (2014).

Results of a large-scale effort of the mouse ENCODE consortium, presenting regulatory and transcript maps of the mouse.

- Gerstein, M. B. et al. Integrative analysis of the *Caenorhabditis elegans* genome by the modENCODE project. *Science* **330**, 1775–1787 (2010).
- The modENCODE Consortium et al. Identification of functional elements and regulatory circuits by *Drosophila* modENCODE. *Science* **330**, 1787–1797 (2010).
- Kudron, M. M. et al. The ModERN Resource: genome-wide binding profiles for hundreds of *Drosophila* and *Caenorhabditis elegans* transcription factors. *Genetics* **208**, 937–949 (2018).
- Gorkin, D. U. et al. An atlas of dynamic chromatin landscapes in mouse fetal development. *Nature* <https://doi.org/10.1038/s41586-020-2093-3> (2020).
- He, P. A. The changing mouse embryo transcriptome at whole tissue and single-cell resolution. *Nature* <https://doi.org/10.1038/s41586-020-2536-x> (2020).
- He, Y. et al. Spatiotemporal DNA methylome dynamics of the developing mouse fetus. *Nature* <https://doi.org/10.1038/s41586-020-2119-x> (2020).
- Cheng, Y. et al. Principles of regulatory information conservation between mouse and human. *Nature* **515**, 371–375 (2014).
- Stefflova, K. et al. Cooperativity and rapid evolution of cobound transcription factors in closely related mammals. *Cell* **154**, 530–540 (2013).
- Keilwagen, J., Posch, S. & Grau, J. Accurate prediction of cell type-specific transcription factor binding. *Genome Biol.* **20**, 9 (2019).
- Tang, F., Lao, K. & Surani, M. A. Development and applications of single-cell transcriptome analysis. *Nat. Methods* **8** (Suppl), S6–S11 (2011).
- Buenrostro, J. D. et al. Single-cell chromatin accessibility reveals principles of regulatory variation. *Nature* **523**, 486–490 (2015).
- Hu, B. C.; HuBMAP Consortium. The human body at cellular resolution: the NIH Human Biomolecular Atlas Program. *Nature* **574**, 187–192 (2019).
- Regev, A. et al. The human cell atlas. *eLife* **6**, e27041 (2017).
- Rhoads, A. & Au, K. F. PacBio sequencing and its applications. *Genomics Proteomics Bioinformatics* **13**, 278–289 (2015).
- Arnold, C. D. et al. Genome-wide quantitative enhancer activity maps identified by STARR-seq. *Science* **339**, 1074–1077 (2013).
- Klein, J. C., Chen, W., Gasperini, M. & Shendure, J. Identifying novel enhancer elements with CRISPR-based screens. *ACS Chem. Biol.* **13**, 326–332 (2018).
- Harrow, J. et al. GENCODE: the reference human genome annotation for The ENCODE Project. *Genome Res.* **22**, 1760–1774 (2012).
- Paudyal, A. et al. The novel mouse mutant, *chuzhoi*, has disruption of Ptk7 protein and exhibits defects in neural tube, heart and lung development and abnormal planar cell polarity in the ear. *BMC Dev. Biol.* **10**, 87 (2010).

Acknowledgements We thank S. Moore, E. Cahill, M. Kellis and J. Li for their assistance, and B. Wold for helpful comments. This work was supported by grants from the NIH: U01HG007019, U01HG007033, U01HG007036, U01HG007037, U41HG006992, U41HG006993, U41HG006994, U41HG006995, U41HG006996, U41HG006997, U41HG006998, U41HG006999, U41HG007000, U41HG007001, U41HG007002, U41HG007003, U41HG007234, U54HG006991, U54HG006997, U54HG006998, U54HG007004, U54HG007005, U54HG007010 and UM1HG009442.

Author contributions The role of the NHGRI Project Management Group in the preparation of this paper was limited to coordination and scientific management of the ENCODE consortium. All other authors contributed to the concepts, writing and/or revisions of this manuscript.

Competing interests B.E.B. declares outside interests in Fulcrum Therapeutics, 1CellBio, HiFiBio, Arsenal Biosciences, Cell Signaling Technologies, BioMillenia, and Nohla Therapeutics. P.F. is a member of the Scientific Advisory Boards of Fabric Genomics, Inc. and Eagle Genomics, Ltd. M.P.S. is cofounder and scientific advisory board member of Personalis, SensOmics, Mirvie, Qbio, January, Filtricine, and Genome Heart. He serves on the scientific advisory board of these companies and Genapsys and Jupiter. Z.W. is a cofounder of Rgenta Therapeutics and she serves on its scientific advisory board. R.M.M. is an advisor to DNAnexus and Decheng Capital, and has outside interests in IMIDomics, Accuragen and ReadCoor, Inc. The authors declare no other competing financial interests.

Additional information

Supplementary information is available for this paper at <https://doi.org/10.1038/s41586-020-2449-8>.

Correspondence and requests for materials should be addressed to M.P.S.

Reprints and permissions information is available at <http://www.nature.com/reprints>.

Publisher's note Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.



Open Access This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons license, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons license, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons license and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this license, visit <http://creativecommons.org/licenses/by/4.0/>.

© The Author(s) 2020

Federico Abascal⁹⁵, Reyes Acosta¹¹, Nicholas J. Addleman¹, Jessika Adrian¹, Veena Afzal⁴⁹, Bronwen Aken¹⁹, Jennifer A. Akiyama⁴⁹, Omar Al Jammal¹⁵, Henry Amrhein⁴³, Stacie M. Anderson⁴⁴, Gregory R. Andrews⁴, Igor Antoshechkin⁴³, Kristin G. Ardlie²⁰, Joel Armstrong¹⁰⁰, Matthew Astley⁹⁵, Budhaditya Banerjee⁷⁷, Amira A. Barkal⁸⁶, If H. A. Barnes¹⁹, Iros Barozzi¹⁴⁹, Daniel Barrell¹⁹, Gemma Barson⁹⁵, Daniel Bates¹¹, Ulugbek K. Baymuradov¹, Cassandra Bazile³⁵, Michael A. Beer^{98,99}, Samantha Beik²⁰, M. A. Bender⁷⁶, Ruth Bennett¹⁹, Louis Philip Benoit Bouvrette^{37,38,39}, Bradley E. Bernstein¹⁶, Andrew Berry¹⁹, Anand Bhaskar¹, Alexandra Bignell¹⁹, Steven M. Blue³⁶, David M. Bodine⁴⁴, Carles Boix²⁶, Nathan Boley¹, Tyler Borrmann⁴, Beatrice Borsari²⁹, Alan P. Boyle^{58,59}, Laurel A. Brandsmeier¹⁸, Alessandra Breschi²⁹, Emery H. Bresnick⁹³, Jason A. Brooks⁴, Michael Buckley¹¹, Christopher B. Burge³⁵, Rachel Byron⁷⁶, Eileen Cahill¹⁵, Lingling Cai⁵, Lulu Cao¹, Mark Cartwright⁹⁶, Rosa G. Castanon⁴⁶, Andres Castillo¹¹, Hassan Chaib¹, Esther T. Chan¹, Daniel R. Chee¹¹, Sora Chee⁹, Hao Chen^{4,6}, Huaming Chen⁴⁸, Jia-Yu Chen⁴¹, Songjie Chen¹, J. Michael Cherry¹, Surya B. Chhetri⁷¹⁸, Jyoti S. Choudhary¹⁰⁷, Jacqueline Chrast¹⁰¹, Dongjun Chung⁹², Declan Clarke⁷, Neal A. L. Cody^{37,38,39}, Candice J. Coppola^{17,18}, Julie Coursen¹⁵, Anthony M. D'Ipollito⁴⁶, Stephen Dalton²⁰, Cassidy Danyko³, Claire Davidson¹⁹, Jose Davila-Velderrain²⁶, Carrie A. Davis³, Job Dekker⁷⁷, Alden Deran¹⁰⁰, Gilberto DeSalvo⁴³, Gloria Despacio-Reyes⁹⁵, Colin N. Dewey⁸⁹, Diane E. Dickel⁴⁹, Morgan Diegel¹¹, Mark Diekhans¹⁰⁰, Vishnu Dileep¹⁰⁰, Bo Ding⁵⁰, Sarah Djebali^{29,30}, Alexander Dobin³, Daniel Dominguez³⁵, Sarah Donaldson¹⁹, Jorg Drenkow³², Timothy R. Drescher¹, Yotam Drier²¹, Michael O. Duff⁴, Douglass Dunn¹¹, Catharine Eastman⁶⁰, Joseph R. Ecker^{48,57}, Matthew D. Edwards²⁶, Nicole El-Alfi⁴⁷, Shaimae I. Elhajjaj⁴, Kerri Elkins³⁶, Andrew Emil⁶¹, Charles B. Epstein²⁰, Rachel C. Evans¹⁸, Iakes Ezkurdia¹⁰², Kaili Fan⁴, Peggy J. Farnham⁶³, Nina Farrell²⁰, Elise A. Feingold¹⁵, Anne-Maud Ferreira¹⁰¹, Katherine Fisher-Aylor⁴³, Stephen Fitzgerald⁹⁵, Paul Flícek¹⁹, Chuan Sheng Foo⁷⁹, Kevin Fortier⁴, Adam Frankish¹⁹, Peter Freese⁴⁰, Shaliu Fu⁵, Xiang-Dong Fu⁴¹, Yu Fu^{4,78}, Yoko Fukuda-Yuzawa⁴⁹, Mariateresa Fulcinitti²², Alister P. W. Funnell¹¹, Idan Gabdank¹, Timur Galeev⁷, Mingshi Gao⁴, Carlos Garcia Giron¹⁹, Tyler H. Garvin⁴⁹, Chelsea Anne Gelboin-Kurkhardt³⁶, Grigoriy Georgopolovs¹¹, Mark B. Gerstein⁷, Belinda M. Giardine¹⁰, David K. Gifford²⁶, David M. Gilbert¹⁰⁸, Daniel A. Gilchrist¹⁵, Shawn Gillespie²¹, Thomas R. Gingeras³, Peng Gong⁶², Alvaro Gonzalez⁹⁶, Jose M. Gonzalez²⁹, Peter Good¹⁷, Alon Goren²⁰, David U. Gorkin^{8,9}, Brenton R. Graveley¹⁴, Michael Gray⁹⁵, Jack F. Greenblatt^{61,71}, Ed Griffiths⁹⁵, Mark T. Groudine⁷⁶, Fabian Grubert¹, Mengting Gu⁷, Roderic Guigo²⁹, Hongbo Guo⁶¹, Yu Guo⁶³, Yuchun Guo²⁶, Gamze Gursoy⁷, Maria Gutierrez-Arcelus⁸⁷, Jessica Halow¹¹, Ross C. Hardison¹⁰, Matthew Hardy¹⁹, Manoj Hariharan⁴⁸, Arif Harmanci⁷, Anne Harrington⁴⁹, Jennifer L. Harrow¹⁰⁶, Tatsunori B. Hashimoto²⁶, Richard D. Hasz¹¹², Meital Hatan²⁰, Eric Haugen¹¹, James E. Hayes⁹⁴, Peng He⁴³, Yupeng He⁴⁸, Nastaran Heidari^{1,64}, David Hendrickson²⁰, Elisabeth F. Heuston⁴⁴, Jason A. Hilton¹, Benjamin C. Hitz¹, Abigail Hochman³⁵, Cory Holgren⁶⁵, Lei Hou²⁶, Shuyu Hou⁵, Yun-Hua E. Hsiao⁹⁷, Shanna Hsu²⁰, Hui Huang⁸, Tim J. Hubbard¹⁰⁵, Jack Huey¹, Timothy R. Hughes^{61,73}, Toby Hunt¹⁹, Sean Ibarrientos¹¹, Robbyn Issner²⁰, Mineo Iwata¹¹, Osagie Izuoglu¹⁹, Tommi Jaakkola²⁶, Nader Jameel⁶⁵, Camden Jansen⁴⁷, Lixia Jiang¹, Peng Jiang^{81,82}, Audra Johnson¹¹, Rory Johnson^{29,33}, Irwin Jungreis^{20,26}, Madhura Kadaba⁶⁵, Maya Kasowski¹, Mary Kasparian⁶⁵, Momoe Kato⁴⁹, Rajinder Kaul^{11,13}, Ruptri Kawli¹, Michael Kay¹⁹, Judith C. Keen¹¹³, Sunduz Keles^{88,89}, Cheryl A. Keller²⁰, David Kelley²⁵, Manolis Kellis^{20,26}, Pouya Kheradpour²⁶, Daniel Sunwook Kim¹, Anthony Kirilusha⁴³, Robert J. Klein⁹⁴, Birgit Knoechel^{22,24}, Samantha Kuan⁹, Michael J. Kulik¹⁰⁰, Sushant Kumar⁷, Anshul Kundaje¹, Tanya Kutuyavin¹¹, Julien Lagarde²⁹, Bryan R. Lajoie⁷⁷, Nicole J. Lambert²⁵, John Lazar¹¹, Ah Young Lee⁸, Donghoon Lee⁷, Elizabeth Lee⁴⁹, Jiri Wook Lee¹, Kristen Lee¹, Christina S. Leslie⁹⁶, Shawn Levy¹⁹, Bin Li⁸, Hairi Li⁴¹, Nan Li⁵⁰, Xiangrui Li⁵, Yang I. Li¹, Ying Li⁵, Yining Li¹, Yue Li²⁶, Jin Lian⁶², Maxwell W. Libbrecht⁸⁰, Shin Lin¹, Ying Lin⁶⁶, Dianbo Liu²⁶, Jason Liu⁷, Peng Liu⁸⁹, Tingting Liu⁵³, X. Shirley Liu^{81,82}, Yan Liu⁵, Yiping Liu²⁶, Maria Long¹⁰, Shaoko Lou⁷, Jane Loveland¹⁹, Aiping Lu⁵, Yuheng Lu²⁶, Eric Lécuyer^{37,38,39}, Lijia Ma⁶⁵, Mark Mackiewicz²⁶, Brandon J. Mannon⁴⁹, Michael Mannetti²¹, Deepa Manthavadi⁹⁵, Georgi K. Marinov⁴³, Fergal J. Martin¹⁹, Eugenio Mattei⁴, Kenneth McCue⁴³, Megan McEown¹⁸, Graham McVicker⁴⁸, Sarah K. Meadows¹⁸, Alex Meissner²⁷, Eric M. Mendenhall^{17,18}, Christopher L. Messer¹⁹, Wouter Meuleman¹¹, Clifford Meyer^{81,82}, Steve Miller⁹⁵, Matthew G. Milton⁶⁵, Tejaswini Mishra¹, Dianna E. Moore¹⁸, Helen M. Moore¹¹⁴, Jill E. Moore⁴, Samuel H. Moore¹⁵, Jennifer Moran⁶⁵, Ali Mortazavi⁴⁷, Jonathan M. Mudge¹⁹, Nikhil Munshi²², Rabi Murad⁴⁷, Richard M. Myers¹⁸, Vivek Nandakumar¹, Preetha Nandi¹⁵, Anil M. Narasimha¹, Aditi K. Narayanan¹, Hannah Naughton¹⁵, Fabio C. P. Navarro⁷, Patrick Navas¹¹, Jirjis Nazarovs⁹⁸, Jemma Nelson¹, Shane Neph¹¹, Fidencio Jun Neri¹¹, Joseph R. Nery⁴⁸, Amy R. Nesmith¹⁸, J. Scott Newberry¹⁶, Kimberly M. Newberry¹⁹, Vu Ngo⁵⁰, Rosy Nguyen¹⁸, Thai B. Nguyen³⁶, Tung Nguyen⁵⁰, Andrew Nishida¹¹, William S. Noble¹², Catherine S. Novak⁴⁹, Eva Maria Novoa²⁶, Briana Nuñez¹⁵, Charles W. O'Donnell²⁶, Sara Olson¹⁴, Kathrina C. Onate¹, Ericka Otterman¹¹, Hakan Ozadam⁷⁷, Michael Pagan¹⁵, Tsultrim Palden²⁵, Xinghua Pan^{62,67,68}, Yongjin Park²⁶, E. Christopher Partridge¹⁶, Benedict Paten¹⁰⁰, Florencia Pauli-Behn¹⁸, Michael J. Pazin¹⁵, Baikang Pei⁷, Len A. Pennacchio^{49,54,56}, Alexander R. Perez⁹⁶, Emily H. Perry¹⁹, Dmitri D. Pervouchine^{29,31}, Nishigandha N. Phalke⁴, Quan Pham⁴⁹, Doug H. Phanfield^{69,70}, Ingrid Plajzer-Frick⁴⁹, Gabriel A. Pratt²⁶, Henry E. Pratt⁴, Sebastian Preissl¹⁸, Jonathan K. Pritchard¹, Yuri Pritykin⁹⁶, Michael J. Purcaro⁴, Qian Qin^{23,84}, Giovanni Quinones-Valdez⁸⁷, Ines Rabano³⁶, Ernest Radovan⁶¹, Anil Raj¹, Nisha Rajagopal⁸⁷, Oren Ram²⁰, Lucia Ramirez⁷, Ricardo N. Ramirez²⁷, Dylan Rausch²¹, Soumya Raychaudhuri⁸⁷, Joseph Raymond²⁰, Rozita Razavi⁷¹, Timothy E. Reddy^{45,46}, Thomas M. Reimond⁴, Bing Ren^{8,9}, Alexandre Reymond¹⁰¹, Alex Reynolds¹¹, Suhn K. Rhie⁶³, John Rinn²⁸, Miguel Rivera²¹, Juan Carlos Rivera-Mulia^{108,109}, Brian Roberts¹⁸, Jose Manuel Rodriguez¹⁰³, Joel Rozowsky⁷, Russell Ryan²¹, Eric Rynes¹¹, Denis N. Salins¹, Richard Sandstrom¹¹, Takayo Sasaki¹⁰⁸, Shashank Sathe³⁶, Daniel Savic⁴², Alexandra Scavell¹³, Jonathan Scheiman²⁹, Christoph Schaffner⁹⁵, Jeffery A. Schloss¹⁵,

Frank W. Schmitges⁷¹, Lei Hoon See³, Anurag Sethi⁷, Manu Setty⁹⁶, Anthony Shafer¹¹, Shuo Shan⁴, Eilon Sharon¹, Quan Shen^{62,72}, Yin Shen^{8,51}, Richard I. Sherwood⁸⁷, Minky Shi¹, Sunyoung Shin⁹⁰, Noam Shores²⁰, Kyle Siebenthal¹¹, Cristina Sisu⁷¹⁰⁴, Teri Slifer⁷, Cricket A. Sloan¹, Anna Smith¹¹⁵, Valentina Snetkova⁴⁹, Michael P. Snyder^{1,2}, Damek V. Spacek¹, Sharanya Srinivasan⁸⁷, Rohith Srivas¹, George Stamatoyannopoulos⁷⁵, John A. Stamatoyannopoulos^{112,123}, Rebecca Stanton³⁶, Dave Steffan⁶⁵, Sandra Stehling-Sun¹¹, J. Seth Strattan¹, Amanda Su³⁵, Balaji Sundararaman³⁶, Marie-Marthe Suerer¹⁹, Tahin Syed²⁶, Matt Szykarek⁶⁵, Forrest Y. Tanaka¹, Danielle Tenen²⁰, Mingxiang Teng⁹⁵, Jeffrey A. Thomas¹¹⁶, Dave Toffey⁶⁵, Michael L. Tress¹⁰³, Diane E. Trout⁴³, Gosia Trynka⁹⁵, Junko Tsuji⁴, Sean A. Upchurch⁴³, Oana Ursu¹, Barbara Uszczyńska-Ratajczak^{29,34}, Mia C. Uziel²⁰, Alfonso Valencia¹⁰³, Benjamin Van Biber¹¹, Arjan G. van der Velde^{4,6}, Eric L. Van Nostrand³⁶, Yekaterina Vaydylevich¹⁵, Jesus Vazquez¹⁰², Alec Victorsen⁶⁵, Jost Vielmetter⁴³, Jeff Verstra¹¹, Axel Visel^{49,54,55}, Anna Vlasova²⁹, Christopher M. Vockley²⁰, Simona Volpi¹⁵, Shinyong Vong¹, Hao Wang¹¹, Mengchi Wang⁵⁰, Qin Wang⁵, Ruth Wang³⁶, Tao Wang⁵⁰, Wei Wang⁵⁰, Xiaofeng Wang^{37,38,39}, Yanli Wang⁶³, Nathaniel K. Watson¹, Xintao Wei¹, Zhijie Wei¹⁵, Hendrik Weisser⁹⁵, Sherman M. Weissman⁶², Rene Welch⁸⁹, Robert E. Welikson¹¹, Zhiping Weng^{4,5,6}, Harm-Jan Westra⁷⁷, John W. Whitaker⁵⁰, Collin White¹⁸, Kevin P. White²⁴, Andre Wildberg⁵⁰, Brian A. Williams⁴³, David Wine²⁰, Heather N. Witt⁶³, Barbara Word⁴³, Maxim Wolf⁶, James Wright⁹⁵, Rui Xiao⁴¹, Xinsu Xiao⁹⁷, Jie Xu⁶³, Jinrui Xu⁷, Koon-Kiu Yan⁴⁷, Yongqi Yan¹¹, Hongbo Yang⁶³, Xinqiong Yang¹, Yi-Wen Yang⁹⁷, Galip Gürkan Yardimci¹², Brian A. Yee³⁶, Gene W. Yeo³⁶, Taylor Young⁴, Tianxiong Yu⁵, Feng Yue^{52,53}, Chris Zaleski³, Chongzhi Zang^{81,82,83}, Haoyang Zeng²⁶, Weihua Zeng⁴⁷, Daniel R. Zerbino¹⁹, Jie Zhai¹, Lijun Zhan¹⁴, Ye Zhan⁷⁷, Bo Zhang⁵³, Jialing Zhang⁶², Jing Zhang⁷, Kai Zhang⁵⁰, Lijun Zhang⁵³, Peng Zhang⁵, Qi Zhang⁹¹, Xiao-Ou Zhang⁴, Yanxiao Zhang⁸, Zhizhuo Zhang²⁶, Yuan Zhao⁸, Ye Zheng⁸⁸, Guoqing Zhong⁶¹, Xiao-Qiao Zhou¹⁵, Yun Zhu⁵⁰ & Jared Zimmerman¹⁰⁸

²⁰The Broad Institute of Harvard and MIT, Cambridge, MA, USA. ²¹MGH, Boston, MA, USA.

²²Dana-Farber Cancer Institute, Boston, MA, USA. ²³Harvard Medical School, Boston, MA, USA.

²⁴Boston Children's Hospital, Boston, MA, USA. ²⁵Harvard University, Cambridge, MA, USA.

²⁶Computer Science and Artificial Intelligence Laboratory, Massachusetts Institute of Technology, Cambridge, MA, USA. ²⁷Max Planck Institute for Molecular Genetics, Department of Genome Regulation, Berlin, Germany. ²⁸University of Colorado Boulder, Boulder, CO, USA.

²⁹Bioinformatics and Genomics Program, Centre for Genomic Regulation (CRG), The Barcelona Institute of Science and Technology and Universitat Pompeu Fabra, Barcelona, Spain. ³⁰IRSD, Université de Toulouse, INSERM, INRA, ENVT, UPS, U1220, CHU Purpan, CS60039, Toulouse, France. ³¹Skolkovo Institute for Science and Technology, Moscow, Russia. ³²Functional Genomics, Cold Spring Harbor Laboratory, Woodbury, NY, USA.

³³Department of Clinical Research, University of Bern, Bern, Switzerland. ³⁴International Institute of Molecular and Cell Biology, Warsaw, Poland. ³⁵Department of Biology, Massachusetts Institute of Technology, Cambridge, MA, USA. ³⁶Department of Cellular and Molecular Medicine, Institute for Genomic Medicine, Stem Cell Program, Sanford Consortium for Regenerative Medicine, University of California, San Diego, La Jolla, CA, USA.

³⁷Département de Biochimie et Médecine Moléculaire, Université de Montréal, Montréal, Québec, Canada. ³⁸Division of Experimental Medicine, McGill University, Québec, Canada.

³⁹Institut de Recherches Cliniques de Montréal (IRCM), Montréal, Québec, Canada. ⁴⁰Program in Computational and Systems Biology, Massachusetts Institute of Technology, Cambridge, MA, USA. ⁴¹Department of Cellular and Molecular Medicine, Institute of Genomic Medicine, University of California, San Diego, San Diego, CA, USA. ⁴²Department of Pharmaceutical Sciences, St. Jude Children's Research Hospital, Memphis, TN, USA. ⁴³Division of Biology, California Institute of Technology, Pasadena, CA, USA. ⁴⁴National Human Genome Research Institute, National Institutes of Health, Bethesda, MD, USA. ⁴⁵Department of Biostatistics and Bioinformatics, Duke University, Durham, NC, USA. ⁴⁶Center for Genomic and Computational Biology, Duke University, Durham, NC, USA. ⁴⁷Biological Sciences, University of California, Irvine, Irvine, CA, USA. ⁴⁸Salk Institute for Biological Studies, La Jolla, CA, USA.

⁴⁹Environmental Genomics and Systems Biology Division, Lawrence Berkeley National Laboratory, Berkeley, CA, USA. ⁵⁰Department of Chemistry and Biochemistry, Department of Cellular and Molecular Medicine, UC San Diego, La Jolla, CA, USA. ⁵¹Institute for Human Genetics, Department of Neurology, University of California, San Francisco, San Francisco, CA, USA. ⁵²Department of Biochemistry and Molecular Genetics, Northwestern University Feinberg School of Medicine, Chicago, IL, USA. ⁵³Penn State Health Milton S. Hershey Medical Center, Hershey, PA, USA. ⁵⁴US Department of Energy Joint Genome Institute, Lawrence Berkeley National Laboratory, Berkeley, CA, USA. ⁵⁵School of Natural Sciences, University of California, Merced, Merced, CA, USA. ⁵⁶Comparative Biochemistry Program, University of California, Berkeley, CA, USA. ⁵⁷Howard Hughes Medical Institute, The Salk Institute for Biological Studies, La Jolla, CA, USA. ⁵⁸Department of Computational Medicine and Bioinformatics, University of Michigan, Ann Arbor, MI, USA. ⁵⁹Department of Human Genetics, University of Michigan, Ann Arbor, MI, USA. ⁶⁰Department of Molecular and Cellular Physiology, School of Medicine, Stanford University, Palo Alto, CA, USA. ⁶¹Terrence Donnelly Centre for Cellular and Biomolecular Research, University of Toronto, Toronto, Ontario, Canada. ⁶²Department of Genetics, School of Medicine, Yale University, New Haven, CT, USA.

⁶³Department of Biochemistry and Molecular Medicine, Norris Comprehensive Cancer Center, Keck School of Medicine, University of Southern California, Los Angeles, CA, USA.

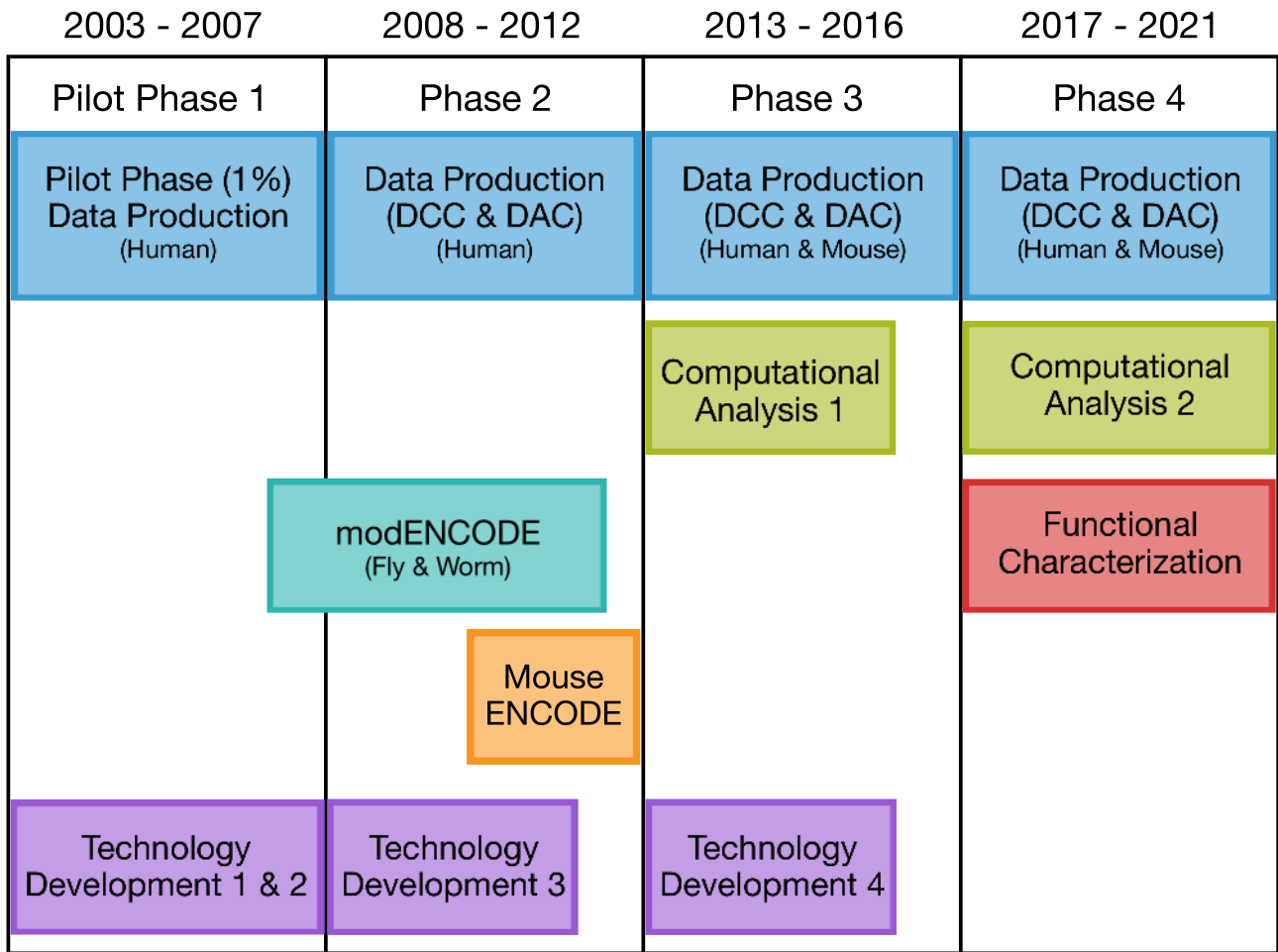
⁶⁴Department of Radiation Oncology, School of Medicine, Stanford University, Palo Alto, CA, USA. ⁶⁵Department of Human Genetics, Institute for Genomics and Systems Biology, The University of Chicago, Chicago, IL, USA. ⁶⁶Division of General Surgery, Section of Transplant Surgery, School of Medicine, Washington University, St. Louis, MO, USA. ⁶⁷Department of Biochemistry and Molecular Biology, School of Basic Medical Sciences, Southern Medical

Perspective

University, Guangzhou, China. ⁶⁸Guangdong Provincial Key Laboratory of Single Cell Technology and Application, Guangzhou, China. ⁶⁹Department of Cell Biology & Physiology, University of North Carolina at Chapel Hill, Chapel Hill, NC, USA. ⁷⁰Thurston Arthritis Research Center, University of North Carolina at Chapel Hill, Chapel Hill, NC, USA. ⁷¹Department of Molecular Genetics, University of Toronto, Toronto, Ontario, Canada. ⁷²School of Medicine, Jiangsu University, Zhenjiang, China. ⁷³Department of Molecular Genetics, Donnelly Centre, University of Toronto, Toronto, Ontario, Canada. ⁷⁴Tempus Labs, Chicago, IL, USA. ⁷⁵Department of Medicine, University of Washington, Seattle, WA, USA. ⁷⁶Fred Hutchinson Cancer Research Center, Seattle, WA, USA. ⁷⁷HHMI and Program in Systems Biology, University of Massachusetts Medical School, Albert Sherman Center, Worcester, MA, USA. ⁷⁸University of Massachusetts Amherst, Amherst, MA, USA. ⁷⁹Institute for Infocomm Research, Singapore, Singapore. ⁸⁰Simon Fraser University, Burnaby, British Columbia, Canada. ⁸¹Center for Functional Cancer Epigenetics, Dana-Farber Cancer Institute, Boston, MA, USA. ⁸²Department of Data Sciences, Dana-Farber Cancer Institute and Harvard T.H. Chan School of Public Health, Boston, MA, USA. ⁸³Center for Public Health Genomics, University of Virginia, Charlottesville, VA, USA. ⁸⁴Molecular Pathology Unit & Cancer Center, Boston, MA, USA. ⁸⁵Department of Biostatistics and Bioinformatics, Moffitt Cancer Center, Tampa, FL, USA. ⁸⁶Institute for Stem Cell Biology and Regenerative Medicine, Stanford University, Stanford, CA, USA. ⁸⁷Department of Medicine, Brigham and Women's Hospital and Harvard Medical School, Boston, MA, USA. ⁸⁸Department of Statistics, Medical Sciences Center, Madison, WI, USA. ⁸⁹Department of Biostatistics and Medical Informatics, Madison, WI, USA. ⁹⁰Department of Mathematical Sciences, University of Texas at Dallas, Richardson, TX, USA. ⁹¹Department of Statistics, University of Nebraska-Lincoln, Lincoln, NE, USA. ⁹²Department of Biomedical Informatics, The Ohio State University, Columbus, OH, USA. ⁹³Department of Cell and

Regenerative Biology, UW-Madison Blood Research Program, Carbone Cancer Center, University of Wisconsin School of Medicine and Public Health, University of Wisconsin, Madison, WI, USA. ⁹⁴Department of Genetics and Genomic Sciences, Icahn School of Medicine at Mount Sinai, New York, NY, USA. ⁹⁵Wellcome Sanger Institute, Cambridge, UK. ⁹⁶Program in Computational Biology, Memorial Sloan Kettering Cancer Center, New York, NY, USA. ⁹⁷610 Charles E. Young Drive S, Terasaki Life Sciences Building, Room 2000E, Los Angeles, CA, USA. ⁹⁸McKusick-Nathans Institute of Genetic Medicine, Johns Hopkins University, Baltimore, MD, USA. ⁹⁹Department of Biomedical Engineering, Johns Hopkins University, Baltimore, MD, USA. ¹⁰⁰University of California, Santa Cruz, Santa Cruz, CA, USA. ¹⁰¹Center for Integrative Genomics, University of Lausanne, Lausanne, Switzerland. ¹⁰²Centro Nacional de Investigaciones Cardiovasculares (CNIC) and CIBER de Enfermedades Cardiovasculares (CIBERCV), Madrid, Spain. ¹⁰³Spanish National Cancer Research Centre (CNIO), Madrid, Spain. ¹⁰⁴Brunel University London, London, UK. ¹⁰⁵King's College London, Guy's Hospital, London, UK. ¹⁰⁶ELIXIR Hub, Wellcome Genome Campus, Cambridge, UK. ¹⁰⁷Institute of Cancer Research, London, UK. ¹⁰⁸Department of Biological Science, Florida State University, Tallahassee, FL, USA. ¹⁰⁹Department of Biochemistry, Molecular Biology and Biophysics, University of Minnesota Medical School, Minneapolis, MN, USA. ¹¹⁰Center for Vaccines and Immunology University of Georgia, Athens, GA, USA. ¹¹¹Center for Molecular Medicine and Department of Biochemistry and Molecular Biology, University of Georgia, Athens, GA, USA. ¹¹²Gift of Life Donor Program, Philadelphia, PA, USA. ¹¹³American Society for Radiation Oncology, Arlington, VA, USA. ¹¹⁴National Cancer Institute, National Institutes of Health, Bethesda, MD, USA. ¹¹⁵Leidos Biomedical, Inc, Frederick, MD, USA. ¹¹⁶National Disease Research Interchange (NDRI), Philadelphia, PA, USA. ¹¹⁷4407 Puller Drive, Kensington, MD, USA.

ENCODE Timeline



Extended Data Fig. 1 | ENCODE timeline. Pilot phase: September 2003–September 2007; ENCODE 2: September 2007–September 2012; ENCODE 3: September 2012–January 2017; ENCODE 4: February 2017–present; modENCODE: April 2007–April 2012; mouse ENCODE: 2009–2012.