



RefSeq Functional Elements as experimentally assayed nongenic reference standards and functional interactions in human and mouse

Catherine M Farrell, Tamara Goldfarb, Sanjida H Rangwala, et al.

Genome Res. published online December 7, 2021

Access the most recent version at doi:[10.1101/gr.275819.121](https://doi.org/10.1101/gr.275819.121)

P<P Published online December 7, 2021 in advance of the print journal.

Accepted Manuscript Peer-reviewed and accepted for publication but not copyedited or typeset; accepted manuscript is likely to differ from the final, published version.

License This is a work of the US Government.

Email Alerting Service Receive free email alerts when new articles cite this article - sign up in the box at the top right corner of the article or [click here](#).



To subscribe to *Genome Research* go to:
<https://genome.cshlp.org/subscriptions>

Published by Cold Spring Harbor Laboratory Press

1 **RefSeq Functional Elements as experimentally assayed non-genic reference standards and**
2 **functional interactions in human and mouse**

3

4 Catherine M. Farrell, Tamara Goldfarb, Sanjida H. Rangwala, Alexander Astashyn, Olga D. Ermolaeva,
5 Vichet Hem, Kenneth S. Katz, Vamsi K. Kodali, Frank Ludwig, Craig L. Wallin, Kim D. Pruitt, and
6 Terence D. Murphy

7

8 National Center for Biotechnology Information, National Library of Medicine, National Institutes of
9 Health, Bethesda, MD 20894, USA

10

11 Corresponding author:

12

13 Catherine M. Farrell, Ph.D.
14 National Center for Biotechnology Information (NCBI)
15 National Library of Medicine
16 National Institutes of Health
17 Bethesda, MD 20894
18 USA

19 E-mail: farrelca@ncbi.nlm.nih.gov

20

21 Running title: NCBI RefSeq Functional Elements

22

23 Keywords: functional element; genome annotation; non-coding genome; non-genic; gene regulation;
24 epigenomic; database; curation

25 **ABSTRACT:**

26 Eukaryotic genomes contain many non-genic elements that function in gene regulation, chromosome
27 organization, recombination, repair or replication, and mutation of those elements can affect genome
28 function and cause disease. While numerous epigenomic studies provide high coverage of gene regulatory
29 regions, those data are not usually exposed in traditional genome annotation, and can be difficult to access
30 and interpret without field-specific expertise. The National Center for Biotechnology Information (NCBI)
31 therefore provides RefSeq Functional Elements (RefSeqFEs), which represent experimentally validated
32 human and mouse non-genic elements derived from the literature. The curated dataset is comprised of
33 richly annotated sequence records, descriptive records in the NCBI Gene database, reference genome
34 feature annotation, and activity-based interactions between non-genic regions, target genes and each
35 other. The dataset provides succinct functional details and transparent experimental evidence, leverages
36 data from multiple experimental sources, is readily accessible and adaptable, and utilizes a flexible data
37 model. The data have multiple uses for basic functional discovery, bioinformatics studies, genetic variant
38 interpretation, as known positive controls for epigenomic data evaluation, and as reference standards for
39 functional interactions. Comparisons to other gene regulatory datasets show that the RefSeqFE dataset
40 includes a wider range of feature types representing more areas of biology, but it is comparatively smaller
41 and subject to data selection biases. RefSeqFEs thus provide an alternative and complementary resource
42 for experimentally assayed functional elements, with future dataset growth expected.

43

44 Eukaryotic genomes contain many types of functional elements, including conventional protein-coding
45 and non-coding genes, gene regulatory elements, architectural elements, and elements associated with
46 DNA replication, recombination and repair. Among those, conventional genes have received the most
47 attention for representation in major genome annotation resources, e.g., RefSeq (O’Leary et al. 2016),
48 GENCODE (Frankish et al. 2021) and others. Gene products, which include alternatively spliced
49 transcripts and proteins, are abundantly represented in genome annotation databases with a heavy focus
50 on protein-coding regions, which occupy less than 1.5% of the mammalian genome. Moreover, genes are
51 major focal points for the curation of disease-associated genetic variation, where there is an emphasis on
52 anchoring variation and linking human disease to specific genes (Wang et al. 2010; Vihinen et al. 2016;
53 Xin et al. 2016; Rivera-Munoz et al. 2018; Amberger et al. 2019; Landrum et al. 2020). Aside from the
54 obvious need to identify gene products due to their importance in biology, such a gene-centric focus is not
55 surprising given that genes and gene-associated variation are generally more amenable for discovery and
56 experimentation than non-genic functional elements, and they tend to offer more tangible avenues for
57 therapeutic treatment of human disease.

58 The genome includes many non-genic elements that function in diverse biological processes, including
59 gene regulation, chromosome organization, recombination or replication. Genome function can be
60 adversely affected by mutation of those elements and result in disease (Lupianez et al. 2016; Chatterjee
61 and Ahituv 2017; Perenthaler et al. 2019; Nesta et al. 2020), supported by genome-wide association
62 studies (GWASs) showing that >90% of disease-associated variation occurs outside of coding regions
63 (Ward and Kellis 2012; Gusev et al. 2014; Albert and Kruglyak 2015; Visscher et al. 2017; Gallagher and
64 Chen-Plotkin 2018; Boix et al. 2021). While a lot of progress has been made in characterizing non-genic
65 functional elements in specialist research fields, that information is not always adequately disseminated to
66 other research fields, most notably to bioclinical research that relies on genome annotation for personal
67 genomic or disease-associated variant interpretation (Perenthaler et al. 2019). Gene regulatory elements
68 are the most abundantly studied among the non-genic element types and their epigenetic signatures are

69 indicated in several large-scale resources, including the Encyclopedia of DNA Elements (The ENCODE
70 Project Consortium 2012), NIH Roadmap Epigenomics (Roadmap Epigenomics Consortium et al. 2015),
71 International Human Epigenome Consortium (Stunnenberg et al. 2016), Ensembl Regulation (Zerbino et
72 al. 2016) and EpiMap (Boix et al. 2021) projects, among others (Garda et al. 2021). However, those data
73 are not usually exposed in traditional genome annotation, can be difficult to interpret, and have not been
74 reconciled with region-specific experimental data in the literature. Thus, complexities in epigenomic data
75 and its consumption disadvantages for non-field-specific experts indicate a need for more highly visible
76 and easily accessible annotation of the non-coding genome. Furthermore, because genome function can
77 only be truly elucidated by taking the entire genome into account, the current gene-centric focus of
78 traditional genome annotation and variant curation points to a general need for better definitions of non-
79 genic functional regions and an ethos to move beyond the genes.

80 To address this, NCBI created RefSeq Functional Elements (RefSeqFEs), a literature-derived dataset that
81 provides reference genome annotation of experimentally validated and well-characterized non-genic
82 regions in human and mouse. The dataset also links functional regions to target genes and to each other
83 when there is activity-based support for functional interactions. Here we describe the creation of this
84 freely available and readily accessible dataset, its multiple components, access options and uses. We also
85 compare RefSeqFEs to other gene regulatory resources, and we report current dataset statistics with
86 feature and genomic distribution analyses, which provide insights into current dataset content and offer
87 suggestions for future needs.

88 **Results**

89 **Dataset scope and design**

90 To distinguish the non-genic dataset from RefSeq conventional genes, which include protein-coding
91 genes, non-coding genes, pseudogenes and gene segments, we defined RefSeqFEs
92 (<https://www.ncbi.nlm.nih.gov/refseq/functionalelements>; Supplemental Table S1) as any genomic

93 element with experimentally validated function, and which is not otherwise considered a conventional
94 gene. For element types we included gene regulatory elements (e.g., enhancers, protein binding sites),
95 known structural elements (e.g., boundary elements, chromatin conformation-associated regions), and
96 other elements of functional importance (e.g., well-defined recombination hotspots or replication origins).
97 While any experimentally validated non-genic element would fall in scope, including elements from high-
98 throughput experimental studies, we prioritized genomic regions that are implicated in human disease or
99 are otherwise of significant interest to the research community. Since we did not aim to replicate the
100 numerous gene regulatory resources that already exist based on well-processed epigenomic or other
101 multi-omics data (Garda et al. 2021), and because re-processing of available omics-derived data was not
102 feasible for us at this time, we decided that RefSeqFEs, at least in the earlier stages of the project, would
103 be focused on smaller-scale experimental data from the literature. Thus, it would be an alternative but
104 complementary literature-derived resource with an emphasis on functional activity. That approach
105 provides flexibility for representing a wide range of feature types in different areas of biology, fills a void
106 to help reconcile other data resources with traditional experimental data in the literature, and allows for
107 robust functional metadata provision such as direct links to publications. Consequently, the current
108 dataset, which is focused on human and mouse, excludes elements from large-scale epigenomic studies
109 and elements that exist solely based on disease-associated variation. It also excludes elements that have
110 indefinite extents or are very large (tens of kilobase or greater lengths), such as telomeres, centromeres,
111 topologically associating domains (TADs) and their broad boundaries, where those are less tractable for
112 genome annotation.

113 For producing the data, we used the existing platforms and workflows already in place for the RefSeq
114 transcript project to take full advantage of NCBI services, such as NCBI search engines, graphical
115 displays and tools, full indexing and versioning of sequence records, and the ability to update records and
116 genome annotation, including on new genome assemblies. Since all RefSeqs are incorporated in NCBI's
117 Nucleotide database (Benson et al. 2018), RefSeqFE sequences and feature annotations adhere to data

118 standards defined by the International Nucleotide Sequence Database Collaboration (INSDC) (Karsch-
119 Mizrachi et al. 2018) with robust use of those standards and ontologies, including recently introduced
120 controlled vocabularies for ‘regulatory_class’ and ‘recombination_class’ features (Supplemental Table
121 S2A,B), and feature qualifiers for metadata displays. Terms from the Sequence Ontology (SO) (Eilbeck et
122 al. 2005) were additionally used to provide further specificity for features lacking a specific INSDC
123 feature or class, and SO terms were also used to define genome-anchored features in GFF3- and bigBed-
124 formatted download files.

125 The dataset was structured to include the following components: 1) Sequence records with curated
126 underlying feature annotation, represented by genomic RefSeq accessions in NCBI’s Nucleotide database;
127 2) Locus-level curated records to integrate metadata, graphical displays and sequences for the underlying
128 region, represented as biological regions in NCBI’s Gene database (Brown et al. 2015); 3) NCBI genome
129 annotation on the human and mouse reference genome assemblies, represented as annotated features with
130 concise and formatted metadata for download and display; and 4) Interaction data to link biological
131 regions to target genes and each other, represented as pairwise interactions.

132 An overview of the RefSeqFE workflow is shown in Figure 1. Briefly, curation was based on
133 experimental data from the literature, with bulk extraction from external databases for large-scale
134 validated datasets, and with supplemental data from researchers if necessary. RefSeq and Gene database
135 records were curated simultaneously, and those records were used as input for genome annotation by
136 NCBI’s Eukaryotic Annotation Pipeline (Pruitt et al. 2014; McGarvey et al. 2015; Supplemental Table
137 S1, annotation pipeline links). Resulting FTP download files and graphical displays were produced for
138 individual RefSeq sequences, Gene database records and genome-annotated features. Those data were
139 further integrated and linked to NCBI-annotated genes and each other to produce additional FTP
140 download files and displays, including a track hub. The following sections expand upon this workflow
141 and provide more details on the components that make up the dataset, followed by analyses of the dataset
142 and its contents.

143 **Sequence records**

144 Genomic RefSeq sequence records with 'NG_' accession prefixes were created to represent the range of
145 one or more experimentally validated non-genic features. We grouped features that were closely located
146 and functionally related in single RefSeqs, such as multiple adjacent or overlapping regulatory elements.
147 The range of those grouped features was used to define a parental biological region (Supplemental Table
148 S2B), as represented in Gene database records described below. To distinguish these non-genic RefSeqs
149 from other genomic 'NG_' accessions represented by RefSeq, all RefSeqFE accessions are associated
150 with NCBI BioProject accession PRJNA343958 (Barrett et al. 2012) and include the keyword RefSeqFE,
151 as indicated in GenBank flat files (Fig. 2A).

152 Both manual curation and automated or semi-automated methods were used to create the RefSeqs.
153 Functional elements were selected for curation initially based on manually scanning the literature for
154 review articles on element types in scope, e.g., gene regulatory elements, recombination regions or
155 replication origins, then identifying specific well-characterized elements described therein and following
156 links to citations. Additional in-scope elements were identified from targeted searches for elements
157 associated with genes of high biomedical interest (e.g., *ACE2*, *BRCA1*, *CFTR*, *HBB* and other frequently
158 accessed genes in the NCBI Gene database), from searches for publications that employ bulk screening
159 techniques, from specific experimental validation term searches in PubMed or PubMed Central, and
160 through outreach efforts and user requests. The current dataset has some inevitable biases for
161 experimentally validated elements that are easily findable (e.g., have been discussed in reviews or are
162 associated with a biomedically important gene), for readily apparent evidence that is well-described and
163 presented in main text of open access publications with a PubMed ID, and for data that are
164 straightforward to curate directly from the publication. Additional details of the data selection and
165 curation process are provided in Supplemental Material. All data were derived from evidence in the
166 literature, either based on individual locus studies or on experimentally validated subsets from larger-
167 scale studies. Examples of high-throughput evidence types used include but are not limited to clustered

168 regularly interspaced short palindromic repeats interference (CRISPRi) assays (e.g., Fulco et al. 2019;
169 Gasperini et al. 2019), massively parallel reporter assays (MPRAs) (e.g., Kheradpour et al. 2013; Ernst et
170 al. 2016), and reporter or transgenic assays from the VISTA project (Visel et al. 2007), FANTOM5
171 project (Andersson et al. 2014) and other bulk-screened datasets (e.g., Wang et al. 2006; Roh et al. 2007;
172 Petrykowska et al. 2008; Narlikar et al. 2010). Those represent a sampling of the available evidence in
173 scope for curation, where new datasets and many more focused region data are continually being
174 identified and will be added to the RefSeqFE dataset over time.

175 We used a wide range of functional features to represent various element types, as indicated in the feature
176 table on our webpage (www.ncbi.nlm.nih.gov/refseq/functionalelements/#Feature_table) and in
177 Supplemental Table S2A,B. A parental biological region feature was annotated on all RefSeqs in addition
178 to one or more underlying functional features. To standardize the curation process, we used the feature
179 definitions provided by INSDC or SO and established policies for the annotation of each feature type
180 (Supplemental Table S2B, columns E,F). In the vast majority of cases, the annotated feature range was
181 defined by the exact fragment tested in an experimental assay, with a minority of features being based on
182 ranges asserted by authors or by sequence analysis tools (see policies per feature type in Supplemental
183 Table S2B, column F). Overlapping feature annotation was allowed, where each feature with a distinct
184 range or type was treated as a unique entity and annotated separately, thus enabling the end user to see
185 each feature as it was assayed in the linked publication(s). Feature type annotation was strictly based on
186 the activity or characteristics primarily shown by the experimental evidence; for example, a protein
187 binding assay and separate evidence that the bound protein functions in a regulatory activity would be
188 represented by separate but overlapping protein binding and regulatory features, such as the *HBB-LCR*
189 5'HS5 CTCF binding site and overlapping enhancer-blocking element features shown in Figure 2B.
190 Adhering to such guidelines enabled straightforward and consistent annotation decisions among curators.
191 Experimental evidence was displayed in INSDC '/experiment' qualifiers on flat files (Fig. 2B, blue tabs),
192 including an evidence type derived from the Evidence & Conclusion Ontology (ECO) (Giglio et al.

193 2019), followed by relevant publication evidence indicated by PubMed IDs. Other feature qualifiers
194 included ‘/note’ and ‘/function’ for additional descriptive and functional information (e.g., cell type
195 activity details), ‘/db_xref’ to link to the associated Gene database record, and feature type-specific
196 INSDC qualifiers. An example of GenBank flat file feature annotation with qualifier and ontology
197 formatting is shown in Figure 2B.

198 **Gene database records**

199 Whereas the RefSeq records provide standalone annotated sequences with feature-specific metadata
200 stored in various INSDC qualifiers, the Gene database serves as the central location for storing various
201 types of metadata at the locus level, while also integrating sequence, genome annotation and graphical
202 display data. Types of locus-level metadata include nomenclature, the locus type designation, a summary
203 based on a synopsis of information from the literature, related publications, orthology information and
204 other standard Gene database fields, as described previously for conventional genes (Brown et al. 2015).

205 To support the RefSeqFE project we created Gene database records identified by a new ‘biological
206 region’ Gene type (Supplemental Fig. S1, red tab). We also added a new ‘Feature type(s)’ field to indicate
207 the types of underlying features annotated on the associated RefSeq (Supplemental Fig. S1, green tab).
208 Since the provision of official nomenclature by the HUGO Gene Nomenclature Committee (HGNC)
209 (Bruford et al. 2020) or Mouse Genome Informatics (MGI) (Zhu et al. 2015) is generally out of scope for
210 non-genic regions, official nomenclature was included for only a few biological regions that had pre-
211 existing official nomenclature. Otherwise, all names, symbols and descriptions were based on curator
212 derivation from the literature, with default symbols containing a ‘*LOC*’ prefix appended with the integer
213 GeneID assigned to the locus.

214 **Genome annotation**

215 All RefSeqFE features (Supplemental Table S3A,B) were annotated by NCBI’s Eukaryotic Genome
216 Annotation Pipeline together with NCBI’s conventional gene-related features, initially in interim human

217 and mouse annotation releases (ARs) starting in 2017 up to the current ARs on the human GRCh38.p13
218 and mouse GRCm39 reference genome assemblies. Following genome annotation, genomic coordinates
219 for annotated biological regions were propagated to relevant Gene records, both in text and graphical
220 formats. Graphical displays of our genome annotation are described below (Fig. 3). Genome-anchored
221 feature annotation was provided in both GFF3 (Moore et al. 2010) and bigBed (Kent et al. 2010) formats
222 for FTP download (Supplemental Table S1, genome annotation data paths), with further details in the
223 ‘Accessing RefSeq Functional Elements data’ section below.

224 **Interaction data**

225 An important aspect of our non-genic annotations is how these regions interplay with each other and with
226 target genes. Therefore, during data curation we internally tracked regulatory element-to-target gene
227 interactions and recombination partner pairings when there was sufficient experimental support in the
228 literature. For regulatory interactions, our linkages were based on either direct experimental evidence for
229 modulation of target gene promoter activity by methods such as reporter gene assays or transgenesis, or
230 by genetic perturbation assays showing regulatory effects on target gene expression. We excluded
231 linkages based on reporter assays that used heterologous promoters, and based on gene proximity
232 predictions, which may only be accurate less than half of the time (Fulco et al. 2019). Thus, only a fifth of
233 the biological regions are linked to target loci, but these linkages have been experimentally assayed and
234 can be used as reference standards for activity-validated interactions. We also tracked regulatory
235 interactions for biological regions that regulate each other, such as distal enhancer activation of a curated
236 promoter region (e.g., the *CFTR* -44 kb enhancer and the *CFTR* promoter, *LOC111674478* and
237 *LOC111674463*, respectively), or when regulatory elements from distinct biological regions have known
238 cooperative activity (e.g., the *CFTR* -44 kb and +36.6 kb enhancers, *LOC111674478* and *LOC111674479*,
239 respectively). Our recombination partner pairings were based on either experimental evidence for non-
240 allelic homologous recombination (e.g., *LOC106804612* and *LOC106804613* representing alpha-globin
241 recombination regions; Fig. 3B), or on direct assays showing translocations or other reproducible

242 recombination events on both sides of a breakpoint (e.g., the *LOC107980440* and *LOC107963955* major
243 breakpoint regions involved in *BCR-ABL* translocations). Both the gene regulatory and recombination
244 interactions were tracked at the parental biological region level, where they are relevant to at least one but
245 not necessarily all underlying features within the biological region.

246 Following reference genome annotation, we determined genomic coordinates for relevant biological
247 regions and target genes and then assembled the pairwise interactions in bigInteract format (Haeussler et
248 al. 2019), also including a custom column listing supporting publications. These data are available for
249 download on our FTP site (<https://ftp.ncbi.nlm.nih.gov/refseq/FunctionalElements/trackhub/data/>;
250 Supplemental Table S1) and can also be visualized in regulatory interaction and recombination tracks in
251 our track hub (Fig. 3B) described below.

252 **Accessing RefSeq Functional Elements data**

253 We provided a variety of data access options for different levels of our data, including for individual
254 RefSeq and Gene database records, and for further processed genome annotation and interaction data. We
255 employed FAIR (Findable, Accessible, Interoperable, Reusable) data principles (Wilkinson et al. 2016) to
256 incorporate compatibility across multiple NCBI and non-NCBI tools and platforms. Our access options
257 are summarized in Supplemental Table S1, where various links are provided for data downloads, sample
258 queries and relevant help documentation. Options are available to access RefSeqFEs via NCBI's Gene
259 database, Nucleotide database, BLAST searching, the BioProject database, NCBI graphical displays, the
260 RefSeqFE Hub (see below), and the NCBI RefSeq, Gene and Genomes FTP sites. In addition, we
261 periodically announce news about the dataset in the NCBI Insights blog
262 (<https://ncbiinsights.ncbi.nlm.nih.gov/tag/refseq-functional-elements/>) and other NCBI social media.

263 To visualize the non-genic biological regions and features, multiple graphical displays were provided for
264 standalone RefSeqs and their genome-annotated contexts (Fig. 3). Each standalone RefSeqFE record can
265 be viewed in graphical format (Supplemental Fig. S2) via a 'Graphics' link at the top of each flat file

266 (Rangwala et al. 2021). Genome-annotated features are color coded according to feature class, and
267 displayed in a ‘Biological regions, aggregate’ track for the indicated NCBI AR (Fig. 3A). The track can
268 be viewed in NCBI graphical view embeds (e.g., in Gene records) and in NCBI’s Genome Data Viewer
269 (GDV) (www.ncbi.nlm.nih.gov/genome/gdv/; Rangwala et al. 2021), enabling the features to be viewed
270 in the context of other data tracks such as variation data, user-uploaded data, remotely connected files or
271 track hubs.

272 To expand the range of genome browsers RefSeqFE annotations can be viewed in and to graphically
273 display the interaction data, we also created a RefSeqFE track hub (Fig. 3B; Supplemental Material). It is
274 in UCSC track hub format (Raney et al. 2014) and serves as a gateway for data visualization, extraction,
275 download and interoperability. It is hosted from the RefSeq FTP site (connection URL:
276 <https://ftp.ncbi.nlm.nih.gov/refseq/FunctionalElements/trackhub/hub.txt>), registered in the Track Hub
277 Registry (Aken et al. 2017) and is a Public Hub in the UCSC Genome browser. It provides parental
278 biological region and underlying feature tracks with custom metadata in bigBed format, and separate
279 tracks for regulatory and recombination interactions in bigInteract format. Additional details on the
280 RefSeqFE Hub and NCBI graphical displays are described in Supplemental Material and on our webpage.

281 While some of our access options are applicable for data querying and use at the biological region level,
282 the ability to query and extract genome-annotated features is likely to be of higher interest. We therefore
283 provided features for the entire set of NCBI-annotated features (including conventional genes) in GFF3
284 format with RefSeqFE features indicated in ‘source’ column 2, and for standalone RefSeqFE features in
285 bigBed format, as described for the RefSeqFE Hub. Links to all our downloadable data can be found in
286 Supplemental Table S1 and on our webpage, where we also provide feature and metadata extraction
287 examples (www.ncbi.nlm.nih.gov/refseq/functionalelements/#Feat_extraction).

288 **Current dataset statistics and content**

289 To qualitatively and quantitatively assess the RefSeqFE dataset, we performed multiple analyses to assess
290 the depth of curation used to produce the dataset, to determine the distribution of features and their
291 genomic locations relative to conventional genes, to assess the relevance of the dataset to clinically
292 relevant genes, and to compare the dataset to other available gene regulatory datasets.

293 To determine the depth of curation used to produce the dataset, we quantified the number of publications
294 used for feature evidence and assessed the number of features derived from each publication
295 (Supplemental Table S2C). In total, we used 2,219 distinct publications as evidence for human AR
296 109.20201120 and mouse AR 109 features combined. A broad set of publications were used as evidence
297 for just a few features (e.g., 85% of publications were used for 1-4 features alone), while a small set of
298 publications for large-scale studies contributed greater than 50 features each and were used as evidence
299 for almost half of the features in the dataset, indicating that the dataset contains a good balance between
300 large-scale and focused study evidence. We additionally assessed the biological regions with respect to
301 single or multiple feature presence and according to study type derivation (Supplemental Table S2D).
302 Approximately 21-23% of biological regions contained multiple features, while 70-73% of biological
303 regions contained a single feature derived from a large-scale study. In summary, these analyses indicate
304 that the dataset is deeply curated from diverse publications with a mix of large-scale and focused studies,
305 and they attest to the high volume of laborious literature review used to create the dataset.

306 To assess the wide range of functional features represented in the dataset (Supplemental Table S2A,B),
307 we determined genome coverage and feature distributions following human AR 109.20201120
308 (GRCh38.p13 assembly) and mouse AR 109 (GRCm39 assembly). In total, and not including parental
309 biological region features, we annotated 9,862 features representing 4,450 distinct biological regions
310 across 6.1 Mb in human, and 2,271 features representing 889 distinct biological regions across 2.2 Mb in
311 mouse (Fig. 4E). The number of annotations per feature type and other feature statistics are shown in
312 Supplemental Table S2A. To further summarize feature distributions, we grouped features into four main
313 types: By INSDC regulatory class, recombination class (represented for human only), protein binding

314 sites and miscellaneous others, as charted in Figure 4A,C and indicated in Supplemental Table S2A. In
315 both human and mouse 63-65% of features were regulatory class. Enhancers were by far the most
316 common among regulatory class features, which may reflect a preference for performing enhancer assays
317 in the literature, likely because their epigenetic signatures make them easier to identify. Protein binding
318 sites were the second-most common feature type in the dataset, accounting for 14% of human and 30% of
319 mouse features, indicating that protein binding assays are also popular in the literature, likely due to the
320 molecular-level functional insights they provide. We noted an underrepresentation of silencer features
321 (only 1-5% of regulatory class features) in the dataset, but we expect to increase silencer representation in
322 the near future based on evidence from recent bulk screens (e.g., Huang et al. 2019; Doni Jayavelu et al.
323 2020; Pang and Snyder 2020).

324 We determined feature length distributions and other length-related statistics for all features combined,
325 per feature class and for individual feature types (Fig. 4B,D; Supplemental Figs. S3 and S4; Supplemental
326 Table S2A). The average length for all features was 781 bps for human and 1,125 bps for mouse, with
327 recombination class features being generally the longest at 2,590 bps, and protein binding sites being
328 generally the shortest at 29-35 bps (Fig. 4B,D). Feature length variability was also apparent between
329 individual feature types within each feature class (Supplemental Figs. S3 and S4; Supplemental Table
330 S2A), e.g., locus control regions were longer than other regulatory class features.

331 To assess the genomic distribution of RefSeqFE features relative to conventional genes and gene
332 subregions, RefSeqFE features were first overlapped with annotated gene ranges, which include introns.
333 We found that more than half (53-54%) of the features were gene range-overlapping in both human and
334 mouse (Fig. 5A,B; Supplemental Table S3A). We further assessed the gene-overlapping features relative
335 to gene subparts (exons, introns, CDS and UTR), and the intergenic features relative to gene 5'-proximal
336 (2 kb upstream of transcript starts) or gene 5'-distal regions (Fig. 5A,B; Supplemental Table S3). For
337 gene overlaps, 37% of all features overlapped introns and 16% overlapped exons in both human and
338 mouse. The majority of exon-overlapping features were UTR-overlapping (12% and 15% of the human

339 and mouse datasets, respectively), while 4% of human and 1% of mouse features were CDS-overlapping,
340 indicating that protein-coding regions may have non-coding biological functions too, a point that may
341 impact genetic variant interpretation as described previously (Hirsch and Birnbaum 2015; Ahitiv 2016).
342 Of the intergenic features, approximately two-thirds were gene-distal, corresponding to 33% of all human
343 and 27% of all mouse features. For all genomic locations, feature overlap completeness was generally
344 high (>75% of features showed >80% overlap with relevant genomic subregions overall; Fig. 5C,D;
345 Supplemental Table S3A), especially with larger genomic segments (whole gene ranges, introns,
346 intergenic regions) or for shorter feature classes (Supplemental Fig. S5), while shorter genomic segments
347 (exons, CDS, UTR) tended to show more partial RefSeqFE feature overlaps.

348 Among the genes that overlapped RefSeqFE features, 64-69% were protein-coding, 28-34% were long
349 non-coding RNA (lncRNA) genes, and 2-4% were other biotypes (Fig. 5E; Supplemental Table S3B,C).
350 In total RefSeqFE features overlapped 2,455 and 565 distinct human and mouse genes, respectively. We
351 also determined that 45% of the human overlapping genes were in at least one clinically relevant gene
352 dataset (Supplemental Table S3B, column 6 square bracket indications), where 833 genes were
353 represented in the RefSeqGene (RSG) dataset (Pruitt et al. 2014), 197 in the Locus Reference Genomic
354 (LRG) dataset (Dagleish et al. 2010), and 835 were genes used for pathogenic (or likely pathogenic)
355 variant submissions to the ClinVar database (Landrum et al. 2020). Cumulatively, RefSeqFE features
356 overlapped 13% of clinically relevant genes from those gene datasets combined, and further gene
357 overlaps are expected upon future dataset growth. This indicates that alternative biological roles may be
358 relevant when interpreting genetic variation in genes of clinical interest. We additionally quantified
359 clinically relevant genes represented as target genes in RefSeqFE human regulatory interactions
360 (Supplemental Table S4). Of 667 distinct target genes, 388 (58%) were represented in at least one of the
361 RSG, LRG and ClinVar gene lists. This likely reflects a high focus on clinically relevant genes in the
362 literature and/or our prioritization of clinical genes for regulatory annotation provision.

363 To further assess RefSeqFEs, we compared the dataset to other gene regulatory resources.
364 Notwithstanding the availability of numerous gene regulatory resources (Garda et al. 2021), we selected
365 just a sampling of those for comparison, namely ENCODE candidate *cis*-regulatory elements (cCREs;
366 The ENCODE Project Consortium et al. 2020), Ensembl Regulation (Zerbino et al. 2016), FANTOM5
367 enhancers (Andersson et al. 2014), VISTA enhancers (Visel et al. 2007) and dbSUPER super-enhancer
368 (Khan and Zhang, 2016). Compared to literature-derived RefSeqFEs, the other resources had different
369 data derivation (Fig. 6A; Supplemental Table S5A), including from epigenomic signatures (ENCODE
370 cCREs, Ensembl and dbSUPER), CAGE data (FANTOM5 enhancers) and transgenic assays (VISTA
371 enhancers). Those resources represented only one or a few feature types (Fig. 6A; Supplemental Table
372 S5A) compared to the over 40 feature types covering more areas of biology in the RefSeqFE resource
373 (Supplemental Table S2A,B). RefSeqFE feature lengths were generally on a par with those from the other
374 datasets (Fig. 6C,E; Supplemental Fig. S6A,C,D; Supplemental Table S5A) except for dbSUPER
375 features, which were longer overall (Supplemental Fig. S6A). However, dataset size and genome
376 coverage comparisons (Fig. 6A; Supplemental Fig. S6B; Supplemental Table S5A) show that the current
377 RefSeqFE dataset is considerably smaller than the comparative datasets with the exception of VISTA
378 enhancers, thereby indicating the major limitation of the dataset, as expected based on its literature-
379 derived nature.

380 To determine feature-level similarity, we intersected RefSeqFE features with features in the other
381 datasets, either individually with each dataset or with features from the comparative datasets combined
382 (Fig. 6B,D; Supplemental Table S5B-F). When all RefSeqFE features were compared to all features in the
383 other resources, approximately 80% of RefSeqFE features overlapped a feature(s) in at least one of the
384 other datasets, with higher overlap percentages being apparent with the larger resources. As expected
385 based on the considerably smaller RefSeqFE dataset size, with the exception of the more similarly sized
386 VISTA dataset, these overlaps represented very low percentages of features in the comparative datasets
387 (Supplemental Table S5B,C, columns E-G), indicating that a lot more content can be gleaned from those

388 large-scale datasets. Nevertheless, a fifth of RefSeqFE features did not show any overlap with the
389 comparative datasets (Fig. 6B,D; Supplemental Table S5B,C,F) and a further 8-25% of pairwise overlaps
390 were poor ($\leq 10\%$ of the RefSeqFE feature was overlapped, Supplemental Table S5D,E, column L),
391 indicating that the dataset contains novel content not represented in the other resources. The non-
392 overlapping features were distributed across all feature classes (30% regulatory, 19% recombination, 22%
393 protein-binding and 29% other types in human and mouse combined; Supplemental Table S5F). A higher
394 proportion of RefSeqFE regulatory or enhancer features overlapped features in the other datasets. (Fig.
395 6B,D). We noted better overlap with ENCODE cCRE enhancers than Ensembl enhancers, likely because
396 ENCODE data were used to identify screening candidates in most of the large-scale studies used as
397 evidence for RefSeqFE enhancers. Many RefSeqFE enhancers correlated with promoter flanking regions,
398 CTCF binding sites and promoters that are abundantly represented in the Ensembl dataset (pairwise
399 feature overlaps in Supplemental Table S5D,E), and indeed non-equivalent feature type overlaps existed
400 with all the comparative datasets, likely due to differences in cell type activity, the versatility of gene
401 regulatory elements, or dataset derivation and completeness differences. The enhancer-only comparisons
402 also indicated high similarities with VISTA positive enhancers (Fig. 6B,D; Jaccard statistics in
403 Supplemental Table S5B,C), as expected given that VISTA positive enhancers are incorporated in the
404 RefSeqFE dataset and are a major source of RefSeqFE enhancers in mouse.

405 In summary, comparisons to other gene regulatory resources indicate that RefSeqFEs represent an
406 alternative but smaller resource based on more traditional experimental evidence from the literature. The
407 dataset offers a greater variety of nuanced feature types covering additional areas of biology, the features
408 generally overlap well with features in comparative resources, and the dataset includes content not found
409 in the other resources. Importantly, the currently smaller and more selective RefSeqFE resource should be
410 considered complementary to other gene regulatory resources.

411 **Discussion**

412 We described here a new literature-derived dataset that provides annotation of experimentally assayed
413 non-genic functional elements in human and mouse, and which uses a robust data model with rich but
414 succinct metadata, and with accessibility options for a wide range of researchers. The dataset includes
415 non-genic elements with diverse biological functions, ranging from gene regulatory elements, replication
416 origins, genomic instability and recombination regions, to gene regulatory and recombination partner
417 interactions. To our knowledge, this combination of functional element annotation is not available in
418 other comparative non-genic data resources. The dataset is unique from a biocuration perspective, where
419 we maximized use of INSDC feature types and qualifiers to format descriptive and functional information
420 from hundreds of publications, with all formatting being accessible and extractable from both standalone
421 RefSeqs and genome annotation. Our provision of RefSeq accessions for standalone use enables sequence
422 findability through various NCBI avenues, including from the Nucleotide and Gene databases and by
423 BLAST analysis. These are more consumable for focused genomic region studies without needing
424 genome-scale extraction, e.g., for sequence determination for subsequent experimental assays, or for
425 using with small-scale sequence analysis tools in the absence of high-performance computation.

426 Our integrative approach with respect to literature-derived data combines diverse experimental data types
427 with a unified metadata structure, and it eliminates user need for exhaustive searching of the literature, or
428 the need to remap data types between different genome assembly versions. Integrating different evidence
429 types can also result in stronger evidence and better inform on function than individual evidence types
430 alone. Although the literature-derived dataset is not inclusive of all available data sources and additional
431 support can be gained from evidence in larger complementary resources, we have already observed
432 strengthened functional support in some biological regions based on multiple evidence types. Examples
433 include *LOC110121455*, *LOC112997545* and *LOC111501765*, where we were able to determine the
434 element type based on reporter assay evidence and link to target genes based on CRISPRi evidence.
435 Further such evidence type combinations are likely to yield more functional insights as the dataset grows.

436 The dataset has multiple uses, ranging from basic functional discovery, to genetic variant interpretation,
437 to use as experimentally validated reference standards in multiple bioinformatic and epigenomic studies.
438 Furthermore, the activity-supported interactions can be used as reference standards for gene regulatory or
439 recombination interactions (also see discussion on their intended use below). Our multiple data
440 accessibility options allow data usage through visual inspection on genome browsers (e.g., for region-
441 specific comparisons to other datasets of interest), or through computational methods based on feature,
442 sequence or metadata extraction, where we have incorporated compatibility with multiple tools and
443 platforms. For basic research, the detailed experimental metadata can inform a researcher on experimental
444 approaches for further in-depth characterization of features of interest. Both the feature annotations and
445 target gene linking may be particularly useful for assigning function to clinically relevant genetic variants,
446 and the experimentally validated features can be used as positive controls for assessing calls in various
447 bioinformatic and epigenomic studies. We have already noted some applications of the dataset in diverse
448 studies, including use of the RefSeqs as a source of locus control regions in a bioinformatics study
449 (Sharma et al. 2019), use of the feature annotation for determining a DNase I hypersensitive site location
450 and sequence in a focused research study (Uchida et al. 2019), and use of the mouse enhancer and
451 promoter features to validate ChIP-seq calls in an epigenomic study (Roller et al. 2021). RefSeqFEs have
452 recently become one of the gene regulatory data sources for the GeneHancer resource (Fishilevich et al.
453 2017). The biological region records have also been used in other resources such as the GeneCards
454 database (Stelzer et al. 2016), and some variation resources link to the biological regions when there is
455 variant overlap, including the Medical Genomics Japan Variant Database (MGeND) (Kamada et al. 2019)
456 and ClinVar Miner (Henrie et al. 2018). NCBI's dbSNP database (Sherry et al. 2001) includes placements
457 relative to RefSeqFE 'NG_' accessions for some SNP entries (e.g., rs11036238), while NCBI's ClinVar
458 resource includes biological regions in the 'Gene(s)' tab for some variant records (e.g., Variation
459 ID:96742). Some biological regions are also reported loci for ClinVar submissions (e.g.,
460 *LOC111365204*). Reciprocally, a link to overlapping ClinVar variants can be found in the 'Variation'

461 section for most human biological regions in the Gene database (Supplemental Table S1, biological
462 regions with ClinVar variants link).

463 The RefSeqFE interactions may appear akin to interactions observed in 3D genomics studies, such as 3D-
464 FISH or chromosome conformation capture-based assays (3C, Hi-C and similar derivatives; Kempfer and
465 Pombo 2020). However, RefSeqFE interactions primarily provide basic element-to-target information as
466 opposed to informing on higher order genome structure, and they are not intended to be comprehensive.
467 They are derived either from genetic manipulation evidence for element-to-target activity or from
468 genomic rearrangement characterization, as opposed to 3D genomics studies which assess physical
469 contacts that are usually mapped at high density. Nevertheless, high-density 3D data can be difficult to
470 interpret and visualize, usually requiring different visualization displays, data formats and specialized
471 analysis tools, thus some users requiring basic element-to-target information may find the relatively
472 simple RefSeqFE interactions easier to use, with the main limitation being the low numbers of
473 interactions in the current dataset, notwithstanding future expected growth. As is the case for the
474 RefSeqFE feature data, the interactions provide complementary data based on an alternative data model.
475 While the RefSeqFE dataset has accessibility, visibility and other advantages, it should be noted that this
476 is a growing dataset where many regulatory elements from the literature still need to be curated, or many
477 functional elements still need to be experimentally validated. This results in obvious disadvantages with
478 respect to genome coverage, and consequently, the current dataset is less useful for researchers seeking
479 comprehensive genome-wide data, for which we encourage the use of larger-scale complementary
480 datasets. The current RefSeqFE dataset is more useful for seeking region-specific functional information
481 (when present) or as an experimentally assayed subset for comparative evaluation of larger-scale data.
482 Other limitations include the data selection biases indicated earlier, including selectivity for data that are
483 easier to curate or automate, and our focus on regions that have been assayed in the literature. The
484 literature itself may have limitations that affect data representation, such as absent, incomplete or
485 inaccurate details in published methods. We expect all of the above limitations to decrease over time as

486 the dataset grows. Other limitations include caveats for some experimental evidence types used in the
487 dataset (Catarino and Stark 2018; Perenthaler et al. 2019). For instance, *in vitro* experimental approaches
488 may not always mimic *in vivo* conditions, including lack of an endogenous chromatin environment, use of
489 heterologous promoters and absence of adjacent accessory sequences in reporter gene assays, lack of a
490 chromatin context in direct protein binding assays, or ectopic genomic integrations resulting in altered
491 chromatin landscapes in transgenic assays. Nevertheless, our representation of features based on those
492 evidence types catalogs them on the genome, alerts researchers about their existence and could potentially
493 prompt further in-depth characterization by other approaches. We also note that the majority of RefSeqFE
494 features based on those evidence types were originally identified as screening candidates from
495 epigenomic or other indicative data in supporting publications, or overlapping features based on
496 alternative evidence types may be present, which boosts confidence in them. While we aim to convey as
497 much functional information about each non-genic element as possible, we recommend that users
498 critically assess experimental evidence and its context.

499 Future plans for RefSeqFEs include dataset growth and qualitative improvements based on research
500 community needs. We aim for significant growth over the next several years, and are particularly
501 interested in engaging with researchers who have data suitable for inclusion in the RefSeqFE dataset.
502 Incorporation of improved and evolving high-throughput functional assays will contribute to dataset
503 growth, including multiplex assays given their high-confidence nature, e.g., epigenetic or 3D genomics
504 information combined with activity assays such as the ChIP-STARR-seq method (Barakat et al. 2018).
505 We plan to increase representation of currently underrepresented feature types and to diversify the sources
506 of high-throughput evidence in the dataset. We will also explore ways to incorporate high-value subsets
507 of large-scale multi-omics data, for which we welcome research community input. We will continue to
508 review our access options and make improvements where necessary, e.g., backfilling and improving
509 cell/tissue type activity data, which is currently only accessible as free-text in feature qualifiers, by
510 converting it to an extractable format. We will provide additional details on our webpage and periodically

511 announce any dataset improvements in the NCBI Insights blog. All community feedback is welcome
512 either directly by e-mail, by using the ‘Feedback’ button on the RefSeqFE webpage or through the
513 RefSeq user mail interface (<https://www.ncbi.nlm.nih.gov/projects/RefSeq/update.cgi>).

514 As the dataset continues to grow, we hope that our literature-derived annotations will provide further
515 insights into how genes are regulated and how the genome functions, with a goal to inform on
516 mechanisms of human disease. Now that we are in the exciting era of genomics in the 2020s, our dataset
517 fulfills a timely need in moving traditional genome annotation beyond the genes, and in disseminating
518 non-genic functional annotation to mainstream research in a more accessible format.

519 **Methods**

520 **RefSeq Functional Elements dataset creation**

521 An overview of the dataset and criteria used for data representation are described in the Results, on our
522 webpage (www.ncbi.nlm.nih.gov/refseq/functionalelements/) and in Figure 1. Procedures to provide
523 sequence records, Gene database records, genome annotation, interaction data and graphical displays are
524 described in relevant sections of the Results and on our webpage. Further specific details on each are
525 available in Supplemental Material.

526 **Data analyses**

527 Data analyses were based on ‘RefSeqFE’ source features extracted from GFF3 files for human AR
528 109.20201120 and mouse AR 109 (FTP download paths in Supplemental Table S1). Full-length gene,
529 gene subpart (‘exon’, ‘CDS’) and 2 kb 5’-proximal features were also extracted from the same GFF3
530 files. Publication metrics were based on extraction of supporting PubMed IDs from bigBed feature files
531 for the same ARs. Clinically relevant gene list sources are provided in Supplemental Table S1 and in
532 Supplemental Material. Comparative datasets were obtained and processed as described in Supplemental
533 Material. Standard UNIX command line methods were used together with the BEDTools software

534 package (Quinlan 2014) to extract and count features, to determine genome coverage and feature length
535 statistics, to deduce intron, UTR and intergenic feature subsets, to determine publication-to-feature and
536 biological region-to-feature metrics, to convert to BED format, to perform feature intersections, and to
537 obtain statistics for overlapping genes, clinically relevant genes, comparative datasets and regulatory
538 target genes. Further specific details on each are available in Supplemental Material.

539 **Competing interest statement**

540 The authors declare no competing interests.

541 **Acknowledgements**

542 This work was supported by the Intramural Research Program of the National Institutes of Health,
543 National Library of Medicine. We thank Drs. Donna Maglott and James Ostell for insights and support
544 during the initiation of this project, Dr. Axel Visel and colleagues for VISTA enhancer discussions, Drs.
545 Daniel Camerini-Otero, Florencia Pratto and Kevin Brick for meiotic recombination region discussions,
546 the Sequence Ontology and INSDC teams for assisting in feature term provision and definitions,
547 numerous NCBI teams and colleagues for contributing to NCBI tools and supporting workflows, and
548 UCSC Genome Browser staff including Drs. Maximilian Haeussler and Brian Lee for track hub support.
549 We are indebted to the countless research scientists who published the experimental evidence used to
550 create the RefSeqFE dataset, and to those who contributed suggestions, supplementary details and data
551 clarifications.

552 *Author Contributions:* The project was conceived, designed and managed by C.M.F. with assistance from
553 T.D.M. and K.D.P. Data curation was performed by C.M.F., T.G and S.H.R. Database development and
554 support were led by T.D.M. and carried out by A.A., O.D.E., V.H., K.S.K., V.K.K., F.L., T.D.M., K.D.P.
555 and C.L.W. with input from C.M.F. Webpage documentation was provided by C.M.F. with assistance

556 from S.H.R. and T.D.M. Track hub preparation, associated download file provision, user outreach, all
 557 data analyses and manuscript writing were carried out by C.M.F. with assistance from T.D.M.

558 **References**

- 559 Ahituv N. 2016. Exonic enhancers: proceed with caution in exome and genome sequencing studies.
 560 *Genome Med* **8**: 14.
- 561 Aken BL, Achuthan P, Akanni W, Amode MR, Bernsdorff F, Bhai J, Billis K, Carvalho-Silva D,
 562 Cummins C, Clapham P et al. 2017. Ensembl 2017. *Nucleic Acids Res* **45**: D635-D642.
- 563 Albert FW, Kruglyak L. 2015. The role of regulatory variation in complex traits and disease. *Nat Rev*
 564 *Genet* **16**: 197-212.
- 565 Amberger JS, Bocchini CA, Scott AF, Hamosh A. 2019. OMIM.org: leveraging knowledge across
 566 phenotype-gene relationships. *Nucleic Acids Res* **47**: D1038-D1043.
- 567 Andersson R, Gebhard C, Miguel-Escalada I, Hoof I, Bornholdt J, Boyd M, Chen Y, Zhao X, Schmidl C,
 568 Suzuki T et al. 2014. An atlas of active enhancers across human cell types and tissues. *Nature*
 569 **507**: 455-461.
- 570 Barakat TS, Halbritter F, Zhang M, Rendeiro AF, Perenthaler E, Bock C, Chambers I. 2018. Functional
 571 Dissection of the Enhancer Repertoire in Human Embryonic Stem Cells. *Cell Stem Cell* **23**: 276-
 572 288.e8
- 573 Barrett T, Clark K, Gevorgyan R, Gorelenkov V, Gribov E, Karsch-Mizrachi I, Kimelman M, Pruitt KD,
 574 Resenchuk S, Tatusova T et al. 2012. BioProject and BioSample databases at NCBI: facilitating
 575 capture and organization of metadata. *Nucleic Acids Res* **40**: D57-63.
- 576 Benson DA, Cavanaugh M, Clark K, Karsch-Mizrachi I, Ostell J, Pruitt KD, Sayers EW. 2018. GenBank.
 577 *Nucleic Acids Res* **46**: D41-D47.
- 578 Boix CA, James BT, Park YP, Meuleman W, Kellis M. 2021. Regulatory genomic circuitry of human
 579 disease loci by integrative epigenomics. *Nature* **590**: 300-307.
- 580 Brown GR, Hem V, Katz KS, Ovetsky M, Wallin C, Ermolaeva O, Tolstoy I, Tatusova T, Pruitt KD,
 581 Maglott DR et al. 2015. Gene: a gene-centered information resource at NCBI. *Nucleic Acids Res*
 582 **43**: D36-42.
- 583 Bruford EA, Braschi B, Denny P, Jones TEM, Seal RL, Tweedie S. 2020. Guidelines for human gene
 584 nomenclature. *Nat Genet* **52**: 754-758.
- 585 Bulger M, Schubeler D, Bender MA, Hamilton J, Farrell CM, Hardison RC, Groudine M. 2003. A
 586 complex chromatin landscape revealed by patterns of nuclease sensitivity and histone
 587 modification within the mouse beta-globin locus. *Mol Cell Biol* **23**: 5234-5244.
- 588 Catarino RR, Stark A. 2018. Assessing sufficiency and necessity of enhancer activities for gene
 589 expression and the mechanisms of transcription activation. *Genes Dev* **32**: 202-223.
- 590 Chan PK, Wai A, Philipsen S, Tan-Un KC. 2008. 5'HS5 of the human beta-globin locus control region is
 591 dispensable for the formation of the beta-globin active chromatin hub. *PLoS One* **3**: e2134.
- 592 Chatterjee S, Ahituv N. 2017. Gene Regulatory Elements, Major Drivers of Human Disease. *Annu Rev*
 593 *Genomics Hum Genet* **18**: 45-63.
- 594 Dagleish R, Flicek P, Cunningham F, Astashyn A, Tully RE, Proctor G, Chen Y, McLaren WM, Larsson
 595 P, Vaughan BW et al. 2010. Locus Reference Genomic sequences: an improved basis for
 596 describing human DNA variants. *Genome Med* **2**: 24.
- 597 Dhar V, Nandi A, Schildkraut CL, Skoultchi AI. 1990. Erythroid-specific nuclease-hypersensitive sites
 598 flanking the human beta-globin domain. *Mol Cell Biol* **10**: 4324-4333.
- 599 Doni Jayavelu N, Jajodia A, Mishra A, Hawkins RD. 2020. Candidate silencer elements for the human
 600 and mouse genomes. *Nat Commun* **11**: 1061.
- 601 Eilbeck K, Lewis SE, Mungall CJ, Yandell M, Stein L, Durbin R, Ashburner M. 2005. The Sequence

- 602 Ontology: a tool for the unification of genome annotations. *Genome Biol* **6**: R44.
- 603 The ENCODE Project Consortium. 2012. An integrated encyclopedia of DNA elements in the human
604 genome. *Nature* **489**: 57-74.
- 605 The ENCODE Project Consortium, Moore JE, Purcaro MJ, Pratt HE, Epstein CB, Shores N, Adrian J,
606 Kawli T, Davis CA, Dobin A et al. 2020. Expanded encyclopaedias of DNA elements in the
607 human and mouse genomes. *Nature* **583**: 699-710.
- 608 Ernst J, Melnikov A, Zhang X, Wang L, Rogov P, Mikkelsen TS, Kellis M. 2016. Genome-scale high-
609 resolution mapping of activating and repressive nucleotides in regulatory regions. *Nat Biotechnol*
610 **34**: 1180-1190.
- 611 Farrell CM, West AG, Felsenfeld G. 2002. Conserved CTCF insulator elements flank the mouse and
612 human beta-globin loci. *Mol Cell Biol* **22**: 3820-3831.
- 613 Fishilevich S, Nudel R, Rappaport N, Hadar R, Plaschkes I, Iny Stein T, Rosen N, Kohn A, Twik M et al.
614 2017. GeneHancer: genome-wide integration of enhancers and target genes in GeneCards.
615 *Database (Oxford)* **2017**: bax028.
- 616 Frankish A, Diekhans M, Jungreis I, Lagarde J, Loveland JE, Mudge JM, Sisu C, Wright JC, Armstrong
617 J, Barnes I et al. 2021. Gencode 2021. *Nucleic Acids Res* **49**: D916-D923.
- 618 Fulco CP, Nasser J, Jones TR, Munson G, Bergman DT, Subramanian V, Grossman SR, Anyoha R,
619 Doughty BR, Patwardhan TA et al. 2019. Activity-by-contact model of enhancer-promoter
620 regulation from thousands of CRISPR perturbations. *Nat Genet* **51**: 1664-1669.
- 621 Gallagher MD, Chen-Plotkin AS. 2018. The Post-GWAS Era: From Association to Function. *Am J Hum*
622 *Genet* **102**: 717-730.
- 623 Garda S, Schwarz JM, Schuelke M, Leser U, Seelow D. 2021. Public data sources for regulatory genomic
624 features. *Med Genet* **33**: 167-177.
- 625 Gasperini M, Hill AJ, McFaline-Figueroa JL, Martin B, Kim S, Zhang MD, Jackson D, Leith A,
626 Schreiber J, Noble WS et al. 2019. A Genome-wide Framework for Mapping Gene Regulation
627 via Cellular Genetic Screens. *Cell* **176**: 377-390 e319.
- 628 Giglio M, Tauber R, Nadendla S, Munro J, Olley D, Ball S, Mitraka E, Schriml LM, Gaudet P, Hobbs ET
629 et al. 2019. ECO, the Evidence & Conclusion Ontology: community standard for evidence
630 information. *Nucleic Acids Res* **47**: D1186-D1194.
- 631 Gusev A, Lee SH, Trynka G, Finucane H, Vilhjalmsdottir BJ, Xu H, Zang C, Ripke S, Bulik-Sullivan B,
632 Stahl E et al. 2014. Partitioning heritability of regulatory and cell-type-specific variants across 11
633 common diseases. *Am J Hum Genet* **95**: 535-552.
- 634 Haeussler M, Zweig AS, Tyner C, Speir ML, Rosenbloom KR, Raney BJ, Lee CM, Lee BT, Hinrichs AS,
635 Gonzalez JN et al. 2019. The UCSC Genome Browser database: 2019 update. *Nucleic Acids Res*
636 **47**: D853-D858.
- 637 Henrie A, Hemphill SE, Ruiz-Schultz N, Cushman B, DiStefano MT, Azzariti D, Harrison SM, Rehm
638 HL, Eilbeck K. 2018. ClinVar Miner: Demonstrating utility of a Web-based tool for viewing and
639 filtering ClinVar data. *Hum Mutat* **39**: 1051-1060.
- 640 Hirsch N, Birnbaum RY. 2015. Dual Function of DNA Sequences: Protein-Coding Sequences Function
641 as Transcriptional Enhancers. *Perspect Biol Med* **58**: 182-195.
- 642 Huang D, Petrykowska HM, Miller BF, Elnitski L, Ovcharenko I. 2019. Identification of human silencers
643 by correlating cross-tissue epigenetic profiles and gene expression. *Genome Res* **29**: 657-667.
- 644 Kamada M, Nakatsui M, Kojima R, Nohara S, Uchino E, Tanishima S, Sugiyama M, Kosaki K,
645 Tokunaga K, Mizokami M et al. 2019. MGeND: an integrated database for Japanese clinical and
646 genomic information. *Hum Genome Var* **6**: 53.
- 647 Karsch-Mizrachi I, Takagi T, Cochrane G, International Nucleotide Sequence Database C. 2018. The
648 international nucleotide sequence database collaboration. *Nucleic Acids Res* **46**: D48-D51.
- 649 Kempfer R, Pombo A. 2020. Methods for mapping 3D chromosome architecture. *Nat Rev Genet* **21**: 207-
650 226.
- 651 Kent WJ, Zweig AS, Barber G, Hinrichs AS, Karolchik D. 2010. BigWig and BigBed: enabling browsing
652 of large distributed datasets. *Bioinformatics* **26**: 2204-2207.

- 653 Khan A, Zhang X. 2016. dbSUPER: a database of super-enhancers in mouse and human genome. *Nucleic*
654 *Acids Res* **44**: D164-171.
- 655 Kheradpour P, Ernst J, Melnikov A, Rogov P, Wang L, Zhang X, Alston J, Mikkelsen TS, Kellis M.
656 2013. Systematic dissection of regulatory motifs in 2000 predicted human enhancers using a
657 massively parallel reporter assay. *Genome Res* **23**: 800-811.
- 658 Landrum MJ, Chitipiralla S, Brown GR, Chen C, Gu B, Hart J, Hoffman D, Jang W, Kaur K, Liu C et al.
659 2020. ClinVar: improvements to accessing data. *Nucleic Acids Res* **48**: D835-D844.
- 660 Long Q, Bengra C, Li C, Kutlar F, Tuan D. 1998. A long terminal repeat of the human endogenous
661 retrovirus ERV-9 is located in the 5' boundary area of the human beta-globin locus control region.
662 *Genomics* **54**: 542-555.
- 663 Lupianez DG, Spielmann M, Mundlos S. 2016. Breaking TADs: How Alterations of Chromatin Domains
664 Result in Disease. *Trends Genet* **32**: 225-237.
- 665 McGarvey KM, Goldfarb T, Cox E, Farrell CM, Gupta T, Joardar VS, Kodali VK, Murphy MR, O'Leary
666 NA, Pujar S et al. 2015. Mouse genome annotation by the RefSeq project. *Mamm Genome* **26**:
667 379-390.
- 668 Moore B, Fan G, Eilbeck K. 2010. SOBA: sequence ontology bioinformatics analysis. *Nucleic Acids Res*
669 **38**: W161-164.
- 670 Narlikar L, Sakabe NJ, Blanski AA, Arimura FE, Westlund JM, Nobrega MA, Ovcharenko I. 2010.
671 Genome-wide discovery of human heart enhancers. *Genome Res* **20**: 381-392.
- 672 Nesta AV, Tafur D, Beck CR. 2020. Hotspots of Human Mutation. *Trends Genet* doi:
673 10.1016/j.tig.2020.10.003.
- 674 O'Leary NA, Wright MW, Brister JR, Ciufu S, Haddad D, McVeigh R, Rajput B, Robbertse B, Smith-
675 White B, Ako-Adjei D et al. 2016. Reference sequence (RefSeq) database at NCBI: current
676 status, taxonomic expansion, and functional annotation. *Nucleic Acids Res* **44**: D733-745.
- 677 Pang B, Snyder MP. 2020. Systematic identification of silencers in human cells. *Nat Genet* **52**: 254-263.
- 678 Perenthaler E, Yousefi S, Niggel E, Barakat TS. 2019. Beyond the Exome: The Non-coding Genome and
679 Enhancers in Neurodevelopmental Disorders and Malformations of Cortical Development. *Front*
680 *Cell Neurosci* **13**: 352.
- 681 Petrykowska HM, Vockley CM, Elnitski L. 2008. Detection and characterization of silencers and
682 enhancer-blockers in the greater CFTR locus. *Genome Res* **18**: 1238-1246.
- 683 Pruitt KD, Brown GR, Hiatt SM, Thibaud-Nissen F, Astashyn A, Ermolaeva O, Farrell CM, Hart J,
684 Landrum MJ, McGarvey KM et al. 2014. RefSeq: an update on mammalian reference sequences.
685 *Nucleic Acids Res* **42**: D756-763.
- 686 Quinlan AR. 2014. BEDTools: The Swiss-Army Tool for Genome Feature Analysis. *Curr Protoc*
687 *Bioinformatics* **47**: 11.12.1-11.12.34.
- 688 Raney BJ, Dreszer TR, Barber GP, Clawson H, Fujita PA, Wang T, Nguyen N, Paten B, Zweig AS,
689 Karolchik D et al. 2014. Track data hubs enable visualization of user-defined genome-wide
690 annotations on the UCSC Genome Browser. *Bioinformatics* **30**: 1003-1005.
- 691 Rangwala SH, Kuznetsov A, Ananiev V, Asztalos A, Borodin E, Evgeniev V, Joukov V, Lotov V, Pannu
692 R, Rudnev D et al. 2021. Accessing NCBI data using the NCBI Sequence Viewer and Genome
693 Data Viewer (GDV). *Genome Res* **31**: 159-169.
- 694 Rivera-Munoz EA, Milko LV, Harrison SM, Azzariti DR, Kurtz CL, Lee K, Mester JL, Weaver MA,
695 Currey E, Craigen W et al. 2018. ClinGen Variant Curation Expert Panel experiences and
696 standardized processes for disease and gene-level specification of the ACMG/AMP guidelines for
697 sequence variant interpretation. *Hum Mutat* **39**: 1614-1622.
- 698 Roadmap Epigenomics Consortium, Kundaje A, Meuleman W, Ernst J, Bilenky M, Yen A, Heravi-
699 Moussavi A, Kheradpour P, Zhang Z, Wang J et al. 2015. Integrative analysis of 111 reference
700 human epigenomes. *Nature* **518**: 317-330.
- 701 Roh TY, Wei G, Farrell CM, Zhao K. 2007. Genome-wide prediction of conserved and nonconserved
702 enhancers by histone acetylation patterns. *Genome Res* **17**: 74-81.
- 703 Roller M, Stamper E, Villar D, Izuogu O, Martin F, Redmond AM, Ramachandran R, Harewood L,

704 Odom DT, Flicek P. 2021. LINE retrotransposons characterize mammalian tissue-specific and
705 evolutionarily dynamic regulatory regions. *Genome Biol* **22**: 62.

706 Sharma BS, Swain PK, Verma RJ. 2019. A Systematic Bioinformatics Approach to Motif-Based Analysis
707 of Human Locus Control Regions. *J Comput Biol* **26**: 1427-1437.

708 Sherry ST, Ward MH, Kholodov M, Baker J, Phan L, Smigielski EM, Sirotkin K. 2001. dbSNP: the
709 NCBI database of genetic variation. *Nucleic Acids Res* **29**: 308-311.

710 Stelzer G, Rosen N, Plaschkes I, Zimmerman S, Twik M, Fishilevich S, Stein TI, Nudel R, Lieder I,
711 Mazor Y et al. 2016. The GeneCards Suite: From Gene Data Mining to Disease Genome
712 Sequence Analyses. *Curr Protoc Bioinformatics* **54**: 1.30.1-1.30.33.

713 Stunnenberg HG, International Human Epigenome Consortium, Hirst M. 2016. The International Human
714 Epigenome Consortium: A Blueprint for Scientific Collaboration and Discovery. *Cell* **167**: 1145-
715 1149.

716 Tuan D, Solomon W, Li Q, London IM. 1985. The "beta-like-globin" gene domain in human erythroid
717 cells. *Proc Natl Acad Sci U S A* **82**: 6384-6388.

718 Uchida N, Hsieh MM, Raines L, Haro-Mora JJ, Demirci S, Bonifacino AC, Krouse AE, Metzger ME,
719 Donahue RE, Tisdale JF. 2019. Development of a forward-oriented therapeutic lentiviral vector
720 for hemoglobin disorders. *Nat Commun* **10**: 4479.

721 Vihinen M, Hancock JM, Maglott DR, Landrum MJ, Schaafsma GC, Taschner P. 2016. Human Variome
722 Project Quality Assessment Criteria for Variation Databases. *Hum Mutat* **37**: 549-558.

723 Visel A, Minovitsky S, Dubchak I, Pennacchio LA. 2007. VISTA Enhancer Browser--a database of
724 tissue-specific human enhancers. *Nucleic Acids Res* **35**: D88-92.

725 Visscher PM, Wray NR, Zhang Q, Sklar P, McCarthy MI, Brown MA, Yang J. 2017. 10 Years of GWAS
726 Discovery: Biology, Function, and Translation. *Am J Hum Genet* **101**: 5-22.

727 Wai AW, Gillemans N, Raguz-Bolognesi S, Pruzina S, Zafarana G, Meijer D, Philipsen S, Grosveld F.
728 2003. HS5 of the human beta-globin locus control region: a developmental stage-specific border
729 in erythroid cells. *EMBO J* **22**: 4489-4500.

730 Wang H, Zhang Y, Cheng Y, Zhou Y, King DC, Taylor J, Chiaromonte F, Kasturi J, Petrykowska H,
731 Gibb B et al. 2006. Experimental validation of predicted mammalian erythroid *cis*-regulatory
732 modules. *Genome Res* **16**: 1480-1492.

733 Wang K, Li M, Hakonarson H. 2010. ANNOVAR: functional annotation of genetic variants from high-
734 throughput sequencing data. *Nucleic Acids Res* **38**: e164.

735 Ward LD, Kellis M. 2012. Interpreting noncoding genetic variation in complex traits and human disease.
736 *Nat Biotechnol* **30**: 1095-1106.

737 Wilkinson MD, Dumontier M, Aalbersberg IJ, Appleton G, Axton M, Baak A, Blomberg N, Boiten JW,
738 da Silva Santos LB, Bourne PE et al. 2016. The FAIR Guiding Principles for scientific data
739 management and stewardship. *Sci Data* **3**: 160018.

740 Xin J, Mark A, Afrasiabi C, Tsueng G, Juchler M, Gopal N, Stupp GS, Putman TE, Ainscough BJ,
741 Griffith OL et al. 2016. High-performance web services for querying gene and variant annotation.
742 *Genome Biol* **17**: 91.

743 Zerbino DR, Johnson N, Juetteman T, Sheppard D, Wilder SP, Lavidas I, Nuhn M, Perry E, Raffailac-
744 Desfosses Q, Sobral D et al. 2016. Ensembl regulation resources. *Database (Oxford)* **2016**:
745 bav119.

746 Zhu Y, Richardson JE, Hale P, Baldarelli RM, Reed DJ, Recla JM, Sinclair R, Reddy TB, Bult CJ. 2015.
747 A unified gene catalog for the laboratory mouse reference genome. *Mamm Genome* **26**: 295-304.

748
749
750
751
752
753

754 **Figure legends**

755 **Figure 1.** Workflow for RefSeqFE dataset production. Full cylinders represent databases, the half-
 756 cylinder represents the indicated data source, and rectangles represent actions. Relevant links to additional
 757 information and data access are provided in Supplemental Table S1.

758
 759 **Figure 2.** Example of a biological region RefSeqFE flat file. Segments of RefSeq accession
 760 NG_052895.1 representing the beta-globin locus control region (*HBB-LCR*) are shown. (A) Top section
 761 of the flat file with a link to BioProject accession PRJNA343958 and the ‘RefSeqFE’ keyword outlined in
 762 red. (B) Segment of the feature annotation section. Features are displayed for the 5’HS5 DNase I
 763 hypersensitive site (Tuan et al. 1985; Dhar et al. 1990; Wai et al. 2003), a transcriptional *cis*-regulatory
 764 region (Long et al. 1998), a CTCF binding site (Farrell et al. 2002; Bulger et al. 2003; Chan et al. 2008)
 765 and an enhancer-blocking element (Farrell et al. 2002). Features include ‘/experiment’ qualifiers with
 766 experimental evidence from the literature as indicated by ECO strings and IDs and links to publications
 767 (blue tabs), ‘/note’ qualifiers with descriptive information (grey tabs), ‘/function’ qualifiers describing the
 768 function of each feature where applicable (green tabs), and a ‘/bound_moiety’ qualifier for the protein
 769 binding site (red tab). All features include a ‘/db_xref’ qualifier (black tabs) linking to the biological
 770 region record in the Gene database (GeneID:109580095), and an INSDC class qualifier when relevant
 771 (orange tabs).

772
 773 **Figure 3.** Graphical displays of RefSeqFE data. (A) NCBI Genome Data Viewer display of genome-
 774 annotated features at the human opsin locus control region (*OPSN-LCR*, GeneID:107604627 also shown
 775 in Supplemental Fig. S1). Underlying features are aggregated and displayed in the ‘Biological regions,
 776 aggregate’ track (outlined in red). Depending on user track set options or the entry point to GDV, the
 777 track may need to be turned on via the configuration interface, as detailed on our webpage (Supplemental
 778 Table S1, graphical displays link). Features are color-coded according to class or type. Coordinates are
 779 based on positions on the genome sequence. An example of a mouseover-activated pop-up box is shown
 780 (overlaid grey box). These boxes contain descriptive and functional information (orange tab) including
 781 experimental evidence and links to publications, as well as a ‘Links & Tools’ area (blue tab) linking to the
 782 related Gene database record and to sequences and BLAST analyses. (B) RefSeqFE Hub view of parental
 783 biological regions, underlying features, and gene regulatory and recombination partner interactions in the
 784 UCSC Genome Browser. Regulatory interactions are shown between the alpha-globin locus control
 785 region (*HBA-LCR*, GeneID:106144573) and the downstream *HBZ*, *HBA2*, *HBA2* and *HBQ1* genes (blue
 786 curved lines), while the recombination partners track visualizes recombination (green curved line)

787 between two alpha-globin recombination regions (*LOC106804612* and *LOC106804613*). Parental
 788 biological regions are denoted by black rectangles in the biological regions track, while the features track
 789 uses color coding as described for A. Further item-specific metadata, display options and links to related
 790 data and tools can be found within item- and track-specific details pages. Depending on the density of
 791 interactions in a region, appropriate zoom levels or configuration modes may need to be adjusted, or
 792 specific hub settings such as multi-region view can be used for viewing interactions between distally
 793 located regions.

794

795 **Figure 4.** RefSeqFE feature distributions. (A) Categorized feature counts from human AR 109.20201120
 796 on the GRCh38.p13 genome assembly with grouping by feature class. The pale blue labels indicate the
 797 feature counts per category, where categories and a full breakdown of feature types and counts are
 798 available in Supplemental Table S2A. (B) Boxplot showing feature length distributions for all human
 799 features (light grey) and individual feature classes, with coloring as in A. Some outliers (maximum length
 800 141940) are not displayed because the Y-axis was scaled to better visualize the distributions of shorter
 801 features. Length distributions per feature class are provided in Supplemental Figure S3 with customized
 802 scaling for each class. n = 9862, 1357, 1379, 926 and 6200 sample points. Additional statistics including
 803 minimums, maximums, averages and standard deviations from the mean are provided in Supplemental
 804 Table S2A. (C) Categorized feature counts from mouse AR 109 on the GRCm39 genome assembly as
 805 shown for human in A. (D) Boxplot showing feature length distributions for all mouse features (light
 806 grey) and individual feature classes, as described for human in B. n = 2271, 109, 690 and 1472 sample
 807 points. Additional details are provided in Supplemental Figure S4 and Supplemental Table S2A. (E)
 808 Summary table with overall counts of annotated features, biological region loci and genome coverage for
 809 the indicated AR.

810

811 **Figure 5.** Locations of RefSeqFE features relative to genes. (A) Locations of features from human AR
 812 109.20201120 compared to NCBI annotated genes and subparts from the same AR. The horizontal bar
 813 graph shows the overall locations (gene-overlapping or intergenic), while the bar-of-pie chart shows more
 814 detailed locations. Blue tones denote genes and subparts while grey tones denote intergenic regions. The
 815 pale blue labels indicate overlapping feature counts for each location, as shown for called overlaps in
 816 Supplemental Table S3A. (B) Locations of features from mouse AR 109 as shown for human in A. (C)
 817 Violin plot showing completeness of human RefSeqFE feature overlaps (overlap length/RefSeqFE feature
 818 length) at each gene-relative location (blue- and grey-tone coloring as in A) and cumulative results for all
 819 locations (blue-grey distribution at left). n = 25029, 5468, 2084, 4373, 743, 1735, 5235, 1906 and 3485

820 sample points. Supporting statistics (Fisher p-values, Jaccard statistics, degree of overlap minimums,
821 maximums, averages and standard deviations) are provided in Supplemental Table S3A. (D) Violin plot
822 showing completeness of mouse feature overlaps at each gene-relative location as described for human in
823 *C. n* = 5810, 1249, 502, 981, 97, 459, 1237, 578 and 707 sample points. Supporting statistics are provided
824 in Supplemental Table S3A. (E) Biotype statistics for genes that are overlapped by RefSeqFE features.
825 The count columns indicate the number of distinct genes overlapped by one or more features, while the %
826 total columns indicate percentages of the total number of genes (2455 human, 565 mouse) overlapped by
827 RefSeqFE features for each biotype.

828

829 **Figure 6.** Comparison of RefSeqFEs to other gene regulatory datasets. (A) Overview showing data
830 derivation, feature type representation and current sizes of each dataset on the human GRCh38.p13 and
831 mouse GRCm39 reference assemblies. Additional information for each dataset is provided in
832 Supplemental Table S5A. (B) Bar graph showing human AR 109.20201120 RefSeqFE feature
833 intersections with the indicated datasets, where the Y-axis represents the percent of input RefSeqFE
834 features showing overlap. All features in comparative datasets were intersected with either all RefSeqFE
835 features (medium blue bars), RefSeqFE regulatory features (grey bars) or RefSeqFE enhancer features
836 (light blue bars). Enhancer features from each dataset were additionally intersected with RefSeqFE
837 enhancer features (dark blue bars). Full statistics including input and overlapping feature counts, overlap
838 percentages with respect to each dataset, Fisher p-values and Jaccard statistics are provided in
839 Supplemental Table S5B, with raw intersection output, feature lengths and degrees of overlap with
840 respect to each dataset in Supplemental Table S5D. Datasets showing overlap with each RefSeqFE
841 feature are also indicated in Supplemental Table S3B, column G. (C) Boxplot showing feature length
842 distributions for the indicated human datasets. Some outliers and dbSUPER feature lengths (maximum
843 498572) are not displayed because the Y-axis was scaled to better visualize shorter feature distributions;
844 see Supplemental Figure S6A for a 50 kb Y-axis scale with dbSUPER data included. *n* = 9862, 926535,
845 622457, 63285 and 1989 sample points. Additional statistics including minimums, maximums, averages
846 and standard deviations from the mean are provided in Supplemental Table S5A. (D) Bar graph showing
847 mouse AR 109 RefSeqFE feature intersections with the indicated datasets, as described for human in *B*.
848 Supporting details are provided in Supplemental Tables S3C and S5C,E. (E) Boxplot showing feature
849 length distributions for the indicated mouse datasets, as described for human in *C*. *n* = 2271, 343747,
850 364670, 49802 and 1291 sample points. Supporting details are provided in Supplemental Table S5A and
851 Supplemental Figure S6.

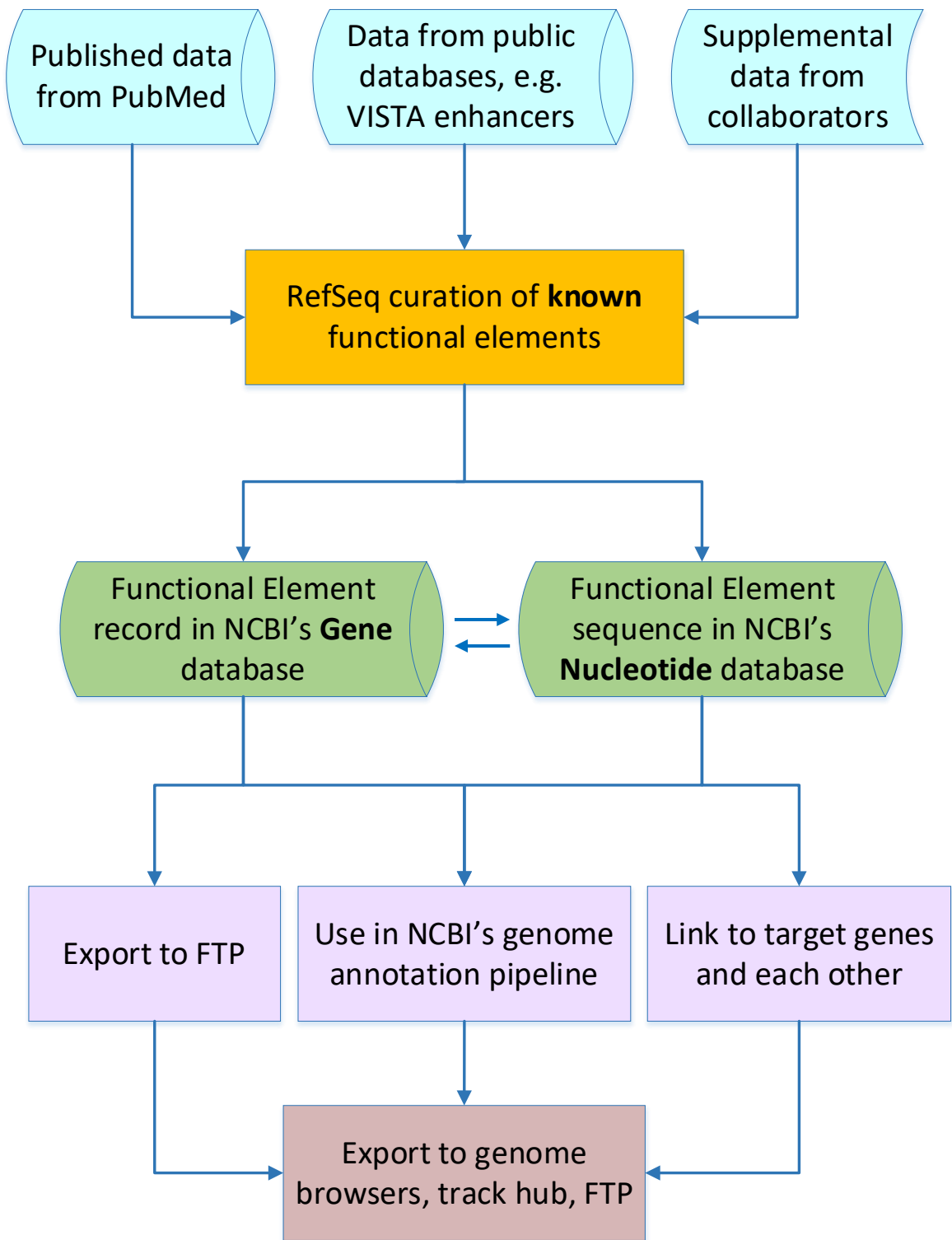


Figure 1

A

LOCUS	NG_052895	34462 bp	DNA	linear	CON 22-APR-2021
DEFINITION	Homo sapiens beta-globin locus control region (HBB-LCR) on chromosome 11.				
ACCESSION	NG_052895				
VERSION	NG_052895.1				
DBLINK	BioProject: PRJNA343958				
KEYWORDS	RefSeq; RefSeqFE.				
SOURCE	Homo sapiens (human)				

B

```

regulatory      20980..21979
                /regulatory_class="DNase_I_hypersensitive_site"
                /experiment="EXISTENCE:in vivo cleavage assay evidence
                [ECO:0001075][PMID:2370867, PMID:3879975, PMID:12941700]"
                /note="5'HS5, also known as HS5, HSS5, HSV or -21.4
                hypersensitive site; not exclusively erythroid; the
                nucleotide coordinates are approximate for this feature"
                /function="enhancer-blocking activity"
                /db_xref="GeneID:109580095"

regulatory      21054..22300
                /regulatory_class="transcriptional_cis_regulatory_region"
                /experiment="EXISTENCE:reporter gene assay evidence
                [ECO:0000049][PMID:9878258]"
                /note="1.2 kb HS5 fragment in the HS5-epsilon-p-CAT
                construct"
                /function="synergizes with the ERV-9 LTR enhancer in
                stably transfected K562 cells"
                /db_xref="GeneID:109580095"

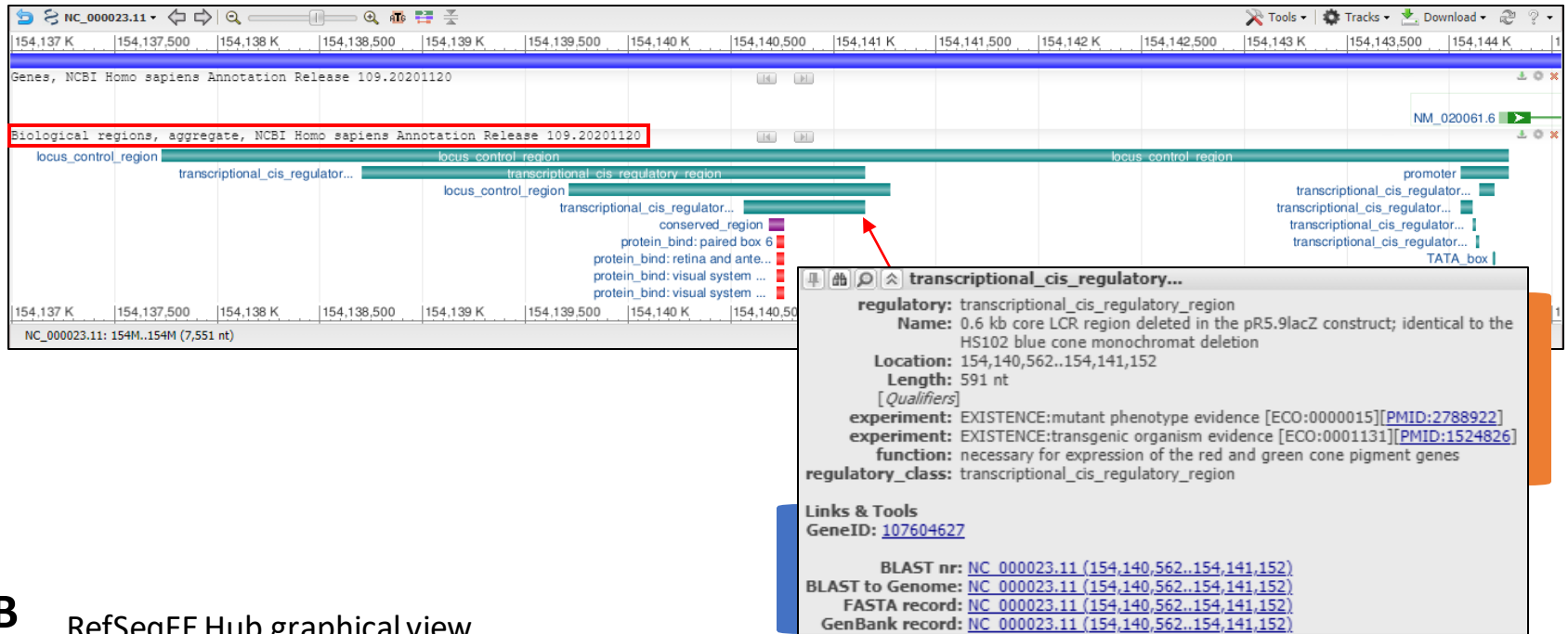
protein bind    21580..21651
                /experiment="EXISTENCE:protein binding evidence
                [ECO:0000024][PMID:11997516, PMID:12861010,
                PMID:18461170]"
                /note="5'HS5 CTCF-binding oligonucleotide"
                /bound_moiety="CCCTC-binding factor"
                /function="enhancer-blocking activity"
                /db_xref="GeneID:109580095"

regulatory      21580..21651
                /regulatory_class="enhancer_blocking_element"
                /experiment="EXISTENCE:reporter gene assay evidence
                [ECO:0000049][PMID:11997516]"
                /note="h5'HS5 enhancer-blocking fragment"
                /function="blocks activation of the Agamma-globin promoter
                by the mouse 5'HS2 enhancer"
                /db_xref="GeneID:109580095"

```

Figure 2

A Genome Data Viewer graphical view



B RefSeqFE Hub graphical view

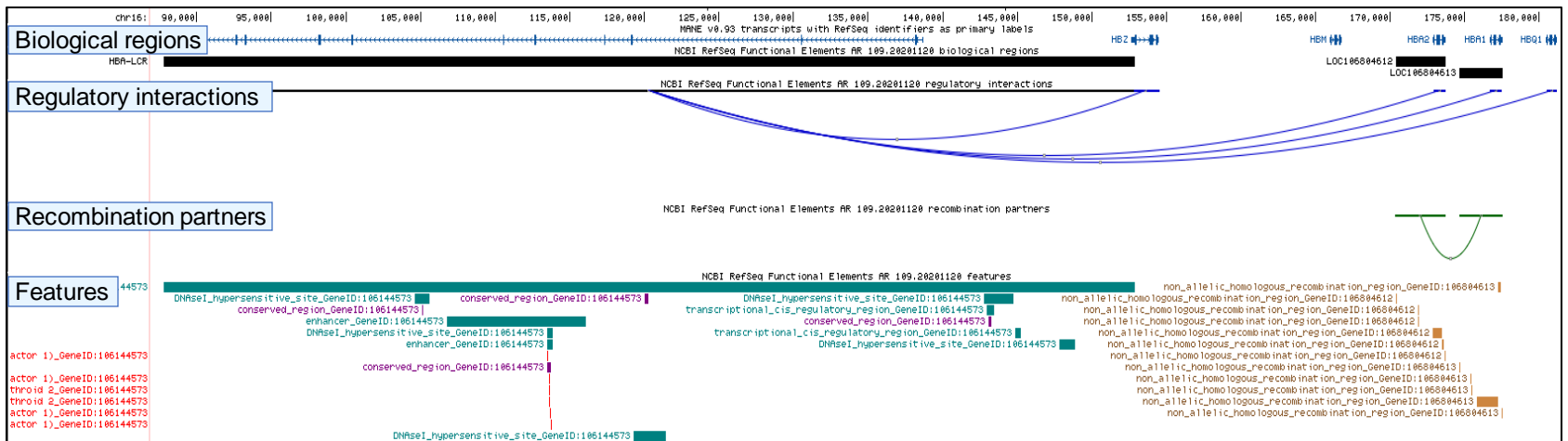


Figure 3

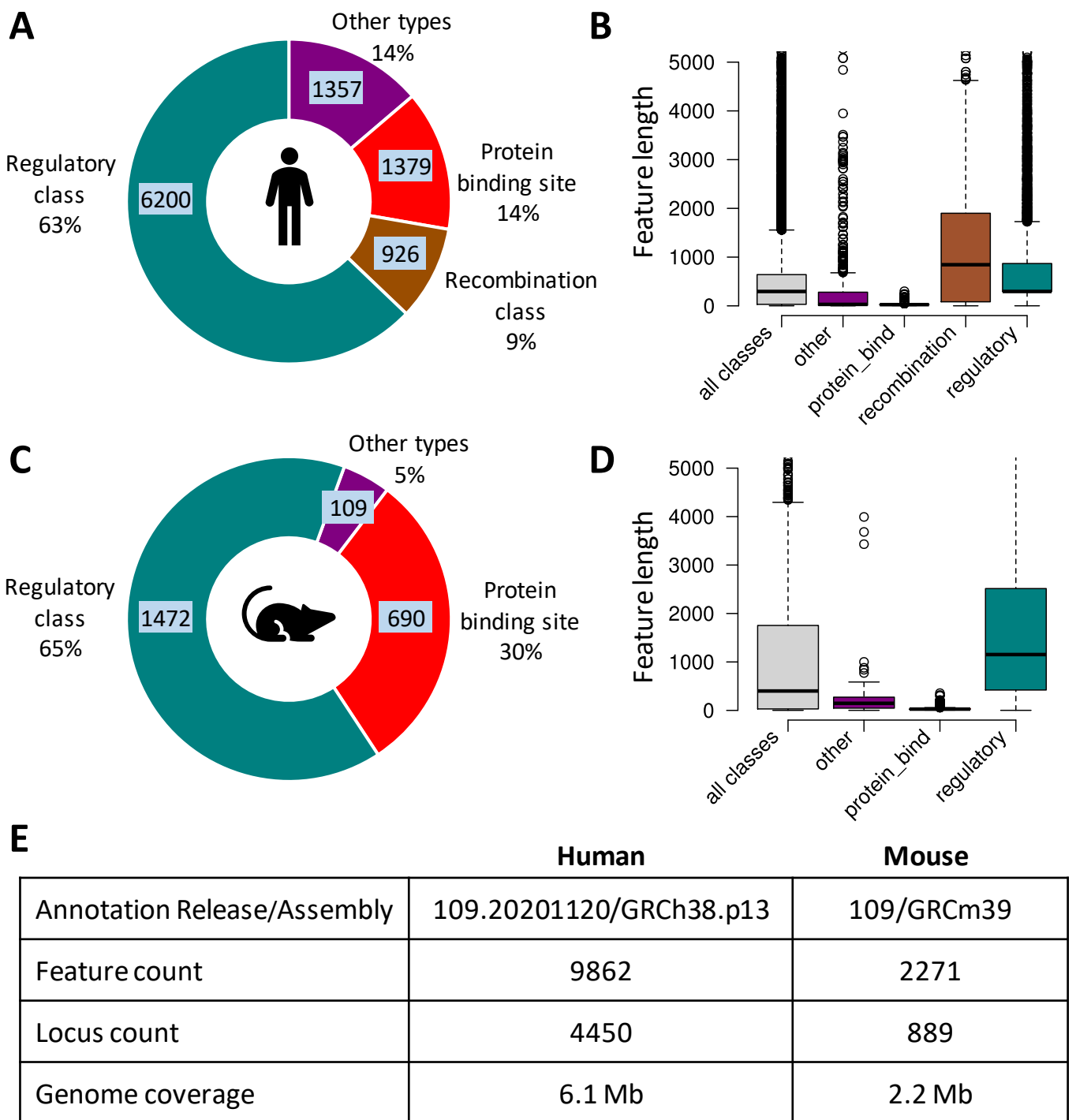


Figure 4

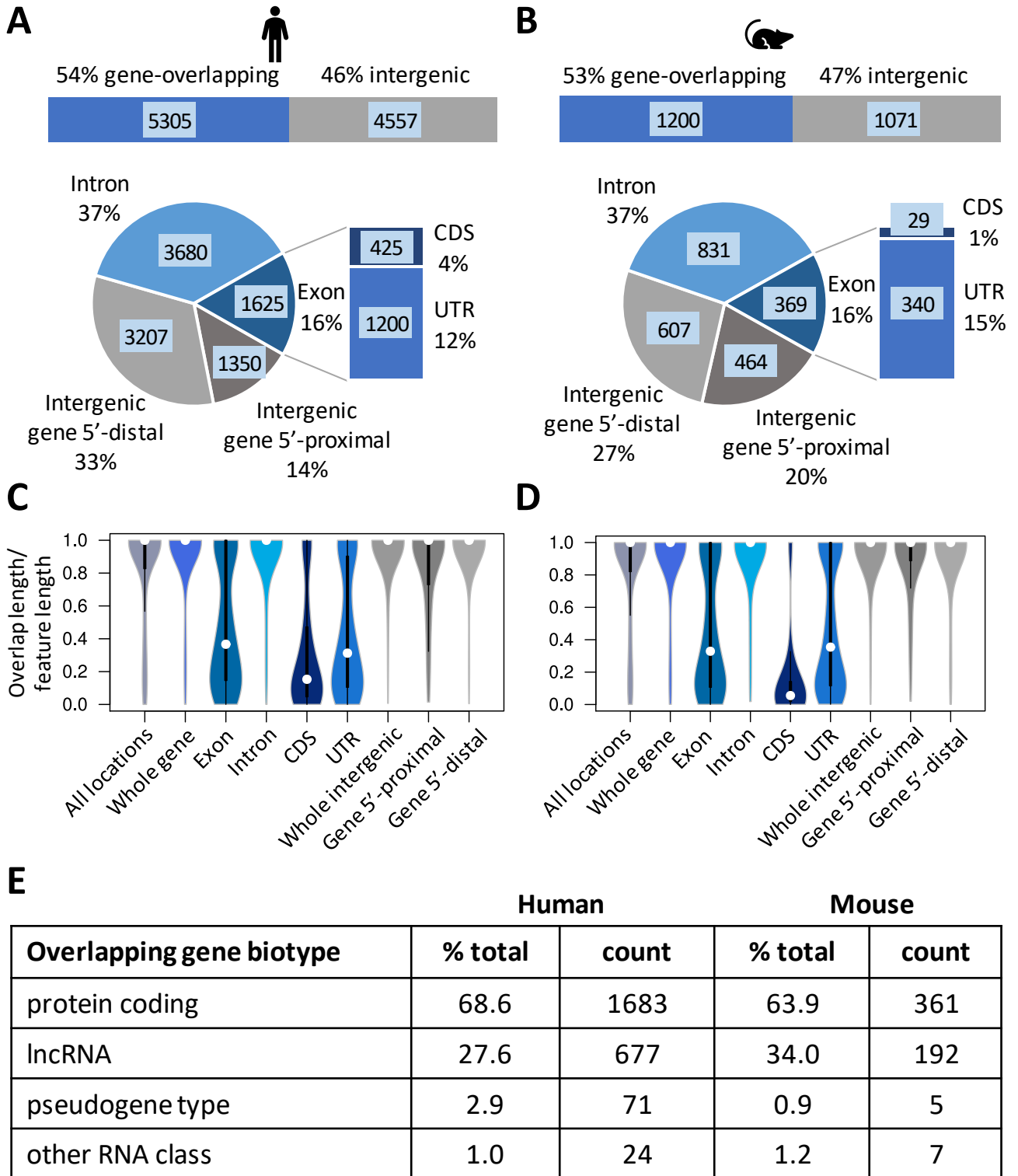
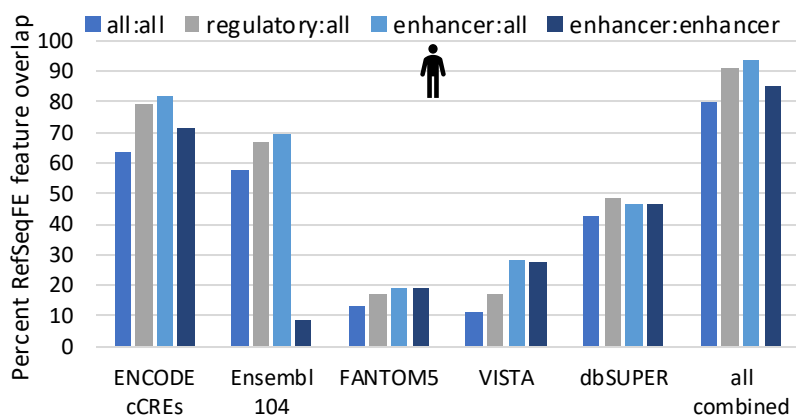
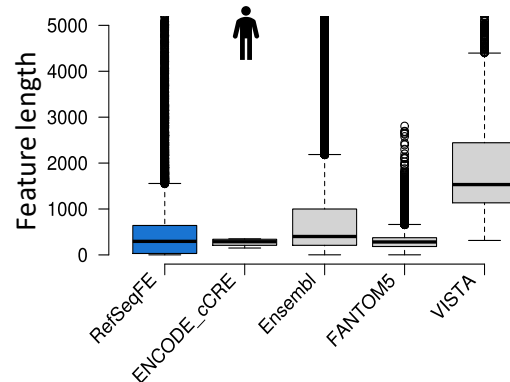
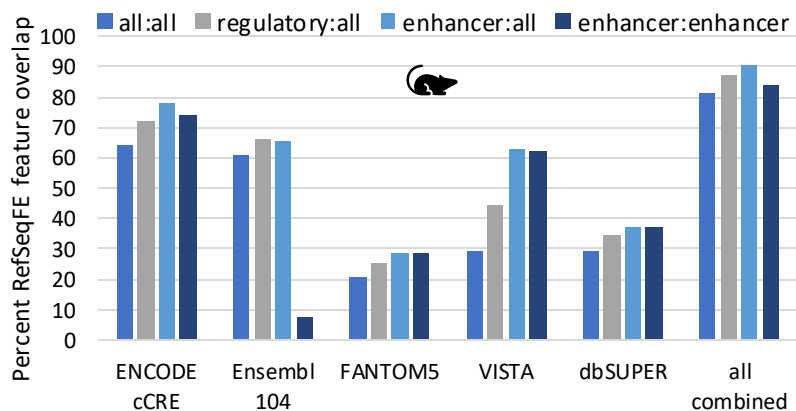
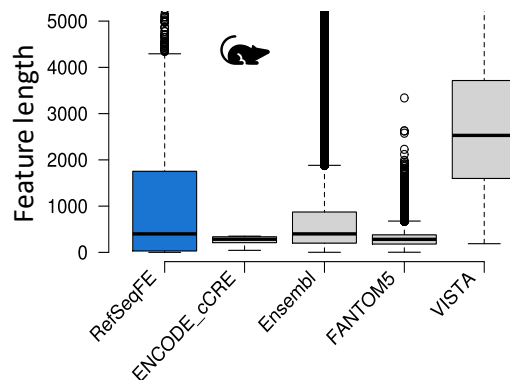


Figure 5

A

Dataset	Data derivation	Feature types	Size
RefSeq Functional Elements	Various experimental assays from the literature	>40 regulatory, protein binding, recombination and other types as in Supplemental Table S2A,B	Human: 9862 features, 6.1 Mb Mouse: 2271 features, 2.2 Mb
ENCODE cCREs	Epigenomic signatures, ChIP-seq, accessibility assays	CTCF-only, DNase-H3K4me3, distal enhancer, promoter, proximal enhancer	Human: 926535 features, 253 Mb Mouse: 343747 features, 298 Mb
Ensembl Regulation	Epigenomic signatures, ChIP-seq, accessibility assays	CTCF binding site, TF binding site, enhancer, open chromatin region, promoter, promoter flanking region	Human: 622457 features, 510 Mb Mouse: 364670 features, 298 Mb
FANTOM5 enhancers	Bidirectional balanced CAGE data, reporter assays for a subset	enhancer	Human: 63285 features, 18.6 Mb Mouse: 49802 features, 14.8 Mb
VISTA enhancers	Comparative genomics, transgenic assays	enhancer	Human: 1989 features, 3.8 Mb Mouse: 1291 features, 3.5 Mb
dbSUPER super-enhancers	Epigenomic signatures, ChIP-seq	super-enhancer	Human: 69340 features, 540 Mb Mouse: 12103 features, 159 Mb

B**C****D****E****Figure 6**