

1 **A sequence-based global map of regulatory activity for deciphering human genetics**

2

3 Kathleen M. Chen^{1,2}, Aaron K. Wong², Olga G. Troyanskaya^{1,2,3*}, Jian Zhou^{4,*}

4

5 ¹Department of Computer Science, Princeton University, Princeton, New Jersey, United States of
6 America

7 ²Flatiron Institute, Simons Foundation, New York, New York, United States of America

8 ³Lewis-Sigler Institute for Integrative Genomics, Princeton University, Princeton, New Jersey,
9 United States of America

10 ⁴Lyda Hill Department of Bioinformatics, University of Texas Southwestern Medical Center,
11 Dallas, Texas, United States of America

12

13 Correspondence to:

14 Olga G. Troyanskaya, email: ogt@cs.princeton.edu

15 Jian Zhou, email: jian.zhou@utsouthwestern.edu

16

17

18

19

20

21

22

23

24

25

26

27

28

29

30

31

32

33

34

35

36

37

38

39

40

41 **Abstract**

42 Sequence is at the basis of how the genome shapes chromatin organization, regulates gene
43 expression, and impacts traits and diseases. Epigenomic profiling efforts have enabled large-
44 scale identification of regulatory elements, yet we still lack a sequence-based map to
45 systematically identify regulatory activities from any sequence, which is necessary for predicting
46 the effects of any variant on these activities. We address this challenge with Sei, a new
47 framework for integrating human genetics data with sequence information to discover the
48 regulatory basis of traits and diseases. Our framework systematically learns a vocabulary for the
49 regulatory activities of sequences, which we call sequence classes, using a new deep learning
50 model that predicts a compendium of 21,907 chromatin profiles across >1,300 cell lines and
51 tissues, the most comprehensive to-date. Sequence classes allow for a global view of sequence
52 and variant effects by quantifying diverse regulatory activities, such as loss or gain of cell-type-
53 specific enhancer function. We show that sequence class predictions are supported by
54 experimental data, including tissue-specific gene expression, expression QTLs, and evolutionary
55 constraints based on population allele frequencies. Finally, we applied our framework to human
56 genetics data. Sequence classes uniquely provide a non-overlapping partitioning of GWAS
57 heritability by tissue-specific regulatory activity categories, which we use to characterize the
58 regulatory architecture of 47 traits and diseases from UK Biobank. Furthermore, the predicted
59 loss or gain of sequence class activities suggest specific mechanistic hypotheses for individual
60 regulatory pathogenic mutations. We provide this framework as a resource to further elucidate
61 the sequence basis of human health and disease.

62

63 **Introduction**

64 Deciphering how regulatory functions are encoded in genomic sequences is a major challenge in
65 understanding how genome variation links to phenotypic traits. Cell-type-specific regulatory
66 activities encoded in elements such as promoters, enhancers, and chromatin insulators are critical
67 to defining the complex expression programs essential for multicellular organisms, like those
68 affecting cell lineage specificity and development. The majority of disease-associated variants
69 from genome-wide association studies (GWAS) are located in noncoding regions¹ and may
70 perturb regulatory elements, yet without knowing how changes in sequence affect regulatory
71 activities we cannot predict the impact of these variants and uncover the regulatory mechanisms
72 contributing to complex diseases and traits. Different variants in the same region can have
73 distinct regulatory consequences and resulting phenotypic effects, as shown by mutations in
74 enhancer regions of SHH²: for instance, a variant may turn off the expression of a gene critical
75 for early development in specific tissue and location, while other variants in the same region may
76 increase enhancer activity or have no effect at all.

77

78 Substantial progress has been made in the experimental profiling and integrative analysis of
79 epigenomic marks, such as histone marks and DNA accessibility, across a wide range of tissues
80 and cell types³⁻⁵. Histone marks are commonly used to identify regulatory elements; for

81 example, H3K4me3 can indicate active promoter regions and H3K27ac/H3K4me1 can indicate
82 active enhancer regions. Moreover, histone marks and chromatin accessibility can be integrated
83 with chromatin state models^{6–10}. These works have been instrumental to annotating the genome
84 with regulatory elements across many tissues.

85
86 At the same time, deep learning sequence modeling techniques have been successfully applied to
87 learn sequence features that are predictive of transcription factor binding and histone
88 modifications^{11–17}. These models are powerful tools for inferring the impact of sequence
89 variation at the chromatin level. However, each chromatin-level prediction can only inform a
90 very specific aspect of sequence--for example, whether a variant causes an increase or decrease
91 of C/EBP- β binding. We continue to lack a global, integrative view of sequence regulatory
92 activities, including all major aspects of cis-regulatory functions, such as tissue-specific or broad
93 enhancer and promoter activities. This limits our ability to interpret the integrated effects of all
94 chromatin-level perturbations caused by genomic variants and determine their impact on human
95 health and diseases.

96
97 We address this challenge by creating a global map for sequence regulatory activity based on a
98 new deep-learning-based framework called Sei. This framework introduces a new sequence
99 model that predicts a comprehensive compendium of 21,907 publicly available chromatin
100 profiles--the broadest set to-date--and uses the model to quantitatively characterize regulatory
101 activities for any sequence with a novel vocabulary of sequence classes. Sequence classes cover
102 diverse types of regulatory activities, such as promoter or cell-type-specific enhancer activity,
103 across the whole genome by integrating sequence-based predictions from histone marks,
104 transcription factors, and chromatin accessibility across a wide range of cell types. For example,
105 ‘embryonic stem cell-specific enhancer’ sequence class activity may be estimated from the
106 predicted binding of multiple transcription factors including Pou5F1, Sox2, and Nanog, as well
107 as various histone marks, on a sequence. Importantly, sequence classes can be used to both
108 classify and quantify the regulatory activities of any sequence based on predictions made by the
109 deep learning sequence model. Therefore, sequence classes allow for the quantitative mapping of
110 any mutation to its impact on cell-type-specific regulatory activities.

111
112 The Sei framework thus provides an interpretable and systematic integration of sequence-based
113 regulatory activity predictions (intrinsic information, based on sequence function) with human
114 genetics data (extrinsic information, based on variant-phenotype association) for discovering the
115 regulatory basis of human traits and disease. We applied our framework to characterize disease-
116 and trait-associated regulatory disruptions by combining sequence class information and UK
117 Biobank GWAS data. Sequence classes provide a non-overlapping partitioning of heritability in
118 GWAS by regulatory activity, which we use to profile the regulatory architecture of 47 diseases
119 and traits in UK Biobank GWAS¹⁸.

120

121 Moreover, variant effect prediction at the sequence-class-level newly enables the interpretation
122 of regulatory mechanisms for individual disease mutations and can differentiate between gain-of-
123 function and loss-of-function regulatory mutations. The regulatory and tissue-specific view
124 provided by sequence classes suggests potential new mechanisms for individual disease-
125 associated variants: for example, we used sequence classes to link mutations in blood-related
126 diseases with previously unknown mechanisms to the malfunctioning of cell-type-specific
127 enhancers.

128
129 We provide the Sei framework as a resource for systematically classifying and scoring any
130 sequence and variant with sequence classes, additionally providing the Sei model predictions for
131 the 21,907 chromatin profiles underlying the sequence classes. The framework can be run using
132 the code at <https://github.com/FunctionLab/sei-framework>, and a user-friendly web server is
133 available at hb.flatironinstitute.org/sei.

134

135 **Results**

136

137 **Developing a comprehensive sequence model for 21,907 chromatin profiles**

138 To capture the widest range of sequence features that are predictive of regulatory activities, we
139 first developed a new deep learning sequence model, which we refer to as the Sei model, that
140 enables the base-level interpretation of sequences by predicting 21,907 genome-wide cis-
141 regulatory targets--including peak calls from 9,471 transcription factor profiles, 10,064 histone
142 mark profiles and 2,372 chromatin accessibility profiles--with single nucleotide sensitivity. The
143 majority of this data (19,905 profiles) is from the Cistrome Project⁵, a resource that uniformly
144 processes and annotates public ChIP-, DNase-, and ATAC-seq datasets, and the remaining
145 chromatin profiles were processed by the ENCODE³ and Roadmap Epigenomics⁴ projects. The
146 Sei model encompasses an estimated ~1000 non-histone DNA-binding proteins (which we refer
147 to as transcription factors), 77 histone marks, and chromatin accessibility across >1300 cell lines
148 and tissues (Supplementary Files 1, 2).

149

150 To efficiently predict 21,907 chromatin profiles from sequence, we designed a novel model
151 architecture (Supplementary Figure 1) and improved our training pipeline. The Sei model uses a
152 new residual-block architecture with a dual linear and nonlinear path design: the linear path allows
153 for fast and statistically efficient training, while the nonlinear path offers strong representation
154 power and the capability to learn complex interactions. For scaling and performance, we
155 introduced a layer of spatial basis functions, which integrates information across spatial locations
156 with much higher memory efficiency than fully connected layers. The model takes as input a 4kb
157 length sequence and predicts the probabilities of 21,907 targets at the center position. The model
158 is trained on chromatin profile peak calls, which are binary (presence/absence), but the model
159 output is continuous, representing probabilities of peaks. Our model training pipeline was updated

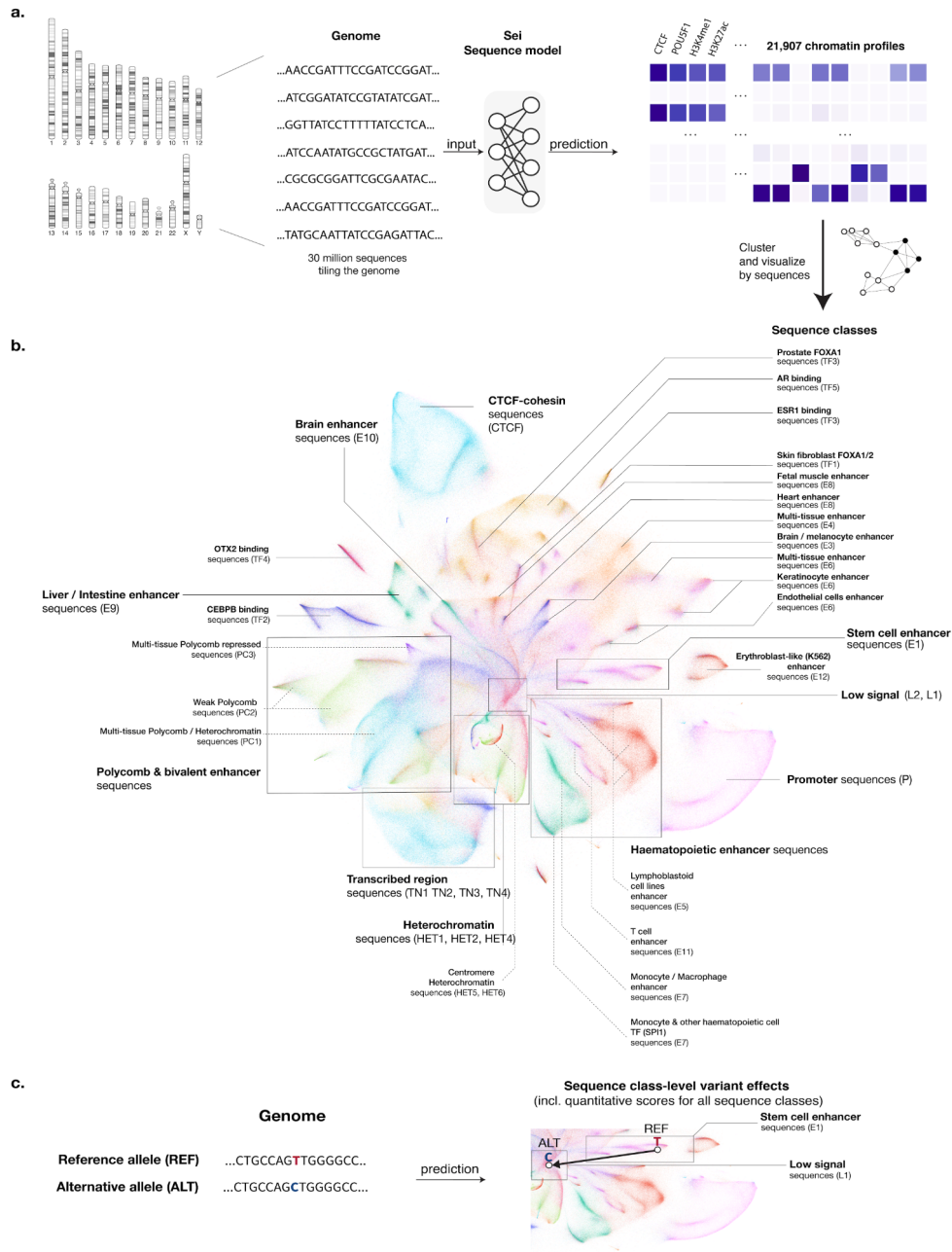
160 to improve training speed and performance by using on-the-fly sampling, which reduces
161 overfitting by generating new training samples for every training step.

162

163 The model achieved an average area under the receiver-operating characteristic (AUROC) of 0.972
164 and average area under the precision-recall curve (AUPRC) of 0.409 across all 21,907 chromatin
165 profiles (Supplementary Figure 2). In addition to accurately predicting individual profiles, the
166 predictions also recapitulated the correlation structure of these profiles, which indicates that the
167 Sei model is able to capture the co-localization patterns of chromatin profiles (Supplementary
168 Figure 3). Furthermore, the Sei model also improved over our best previously published model,
169 DeepSEA “Beluga”¹³, on the 2002 chromatin profiles predicted by both models by 19% on average
170 (as measured by AUROC/1-AUROC, Supplementary Figure 4).

171

172 Therefore, the Sei model is the most comprehensive chromatin-level sequence model to-date, and
173 offers an expansive new resource for sequence and variant interpretation.



174

175 **Figure 1. Mapping the global regulatory landscape of genomic sequences.**

176 **a**, Overview of the Sei framework for systematic prediction of sequence regulatory activities. Sequence
 177 classes are extracted from the predicted chromatin profiles of 30 million sequences evenly tiling the
 178 genome. The predictions were made by Sei, a new deep convolutional network sequence model trained on
 179 21,907 chromatin profiles. Specifically, classes are identified by applying Louvain community detection
 180 to the nearest-neighbor graph of 180 principal components extracted from the predictions data. **b**,
 181 Visualizing the global regulatory landscape of human genome sequences discovered by this approach
 182 with UMAP. Major sequence classes include cell-type-specific enhancer classes, CTCF-cohesin,
 183 promoter, TF-specific, and heterochromatin/centromere classes. **c**, This framework is further applied to
 184 predict sequence-class-level genome variant effects, quantified by changes in sequence class scores.

185

186

187 **Defining sequence classes using a sequence model from whole genome sequences**

188 Next, we applied the Sei model to develop a global, quantitative map from genomic sequences to
189 specific classes of regulatory activities, which we term sequence classes, by integrating the wide
190 range of chromatin profiles predicted by Sei. Sequence classes are therefore mapped directly from
191 sequence, and each sequence class represents a distinct program of regulatory activities across
192 tissues and cell types as covered by the Sei model. Furthermore, sequence classes allow for the
193 mapping of any sequence to quantitative scores that represent a broad spectrum of regulatory
194 activities.

195

196 To cover the whole spectrum of sequence activities, we identified sequence classes from Sei
197 predictions for 30 million sequences uniformly tiling the whole genome (4kb windows with 100bp
198 step size). We visualized the global structure of sequence regulatory signals as represented by the
199 model's chromatin profile predictions with nonlinear dimensionality reduction techniques^{19,20}
200 (Figure 1) and applied Louvain community clustering²¹ to these predictions to categorize the 30
201 million sequences into 40 sequence classes (Figure 1a).

202

203 This visualization of human genome sequences demonstrates the global organization of sequence
204 regulatory activities (Figure 1b). The center of the visualization contains sequences with weak or
205 no regulatory activity based on histone mark and TF enrichment, and sequences with specific
206 regulatory activities radiate outwards, establishing a continuum from no activity to strong specific
207 activity. Different branches of sequences are enriched in distinct chromatin modifications and
208 transcription factors, and sequences with similar regulatory activities are grouped together. For
209 example, tissue-specific enhancer sequences were predominantly grouped by tissue in the
210 visualization (Figure 1b). In addition, sequences with repressive Polycomb marks were spatially
211 adjacent to H3K9me3-marked heterochromatin sequences (Figure 1b), reflecting their extensive
212 crosstalk in epigenetic silencing²²⁻²⁴. Notably, promoter-proximal and CTCF-cohesin binding
213 sequences form two well-defined clusters that are separated from other sequences, which may
214 reflect the distinct nature of these activities (Figure 1b).

215

216 The sequence classes identified from whole genome sequences recapitulate the sequence
217 organization shown in the visualization and provide a basis for summarizing sequence activities
218 globally and are robust to clustering parameter choices (Supplementary Figures 5, 6). To facilitate
219 intuitive interpretation of sequence classes, we named them based on the corresponding
220 enrichment of cis-regulatory profiles (Figure 2a, Supplementary Figures 7-12, Supplementary File
221 3); specifically, we label each sequence class with a functional group acronym and index denoting
222 the rank of the sequence class within the group (Supplementary Figure 13, e.g. E1 encompasses a
223 larger proportion of the genome than E2). Because genomic sequences encode their regulatory
224 activity programs across all cell types, sequence classes also show distinct activity patterns across

225 cell types and tissues. We label sequence classes primarily based on their active, cell-type-specific
226 regulatory activities--in particular, promoter and enhancer activities. Therefore, sequence classes
227 that are not labeled as enhancer ('E') or promoter ('P') generally lack enhancer or promoter activity
228 in any cell type predicted by Sei.

229
230 In summary, sequence classes contain 1 'P' promoter class, which is most strongly enriched in the
231 active promoter histone mark H3K4me3 across all cell types (Figure 2a, Supplementary Figure 7);
232 12 'E' enhancer classes, which are strongly enriched in enhancer histone marks, such as H3K4me1
233 and H3K27ac, and transcription factors relevant to their activities in select cell types (e.g.
234 PU.1/Spi1 in the E7 monocyte/macrophage enhancer class, HNF4- α in E9 liver/intestine, and
235 Sox2/Nanog/Pou5f1 in E1 stem cell), and often display repressive H3K27me3 marks in inactive
236 cell types (Figure 2a, Supplementary Figures 8-10, Supplementary File 3); 1 'CTCF' sequence
237 class, which is strongly enriched in CTCF and cohesin (Figure 2a, Supplementary File 3); 5 'TF'
238 sequence classes, which are enriched in a few specific transcription factors (e.g. CEBPB sequence
239 class) but have weak or no enhancer mark enrichment (Figure 2a, Supplementary File 3); 4 'PC'
240 Polycomb classes, which are enriched in the Polycomb-repressed region mark H3K27me3 and
241 generally not enriched in active promoter or enhancer marks (Figure 2a, Supplementary Figure
242 10); 6 'HET' heterochromatin classes, which are enriched in the heterochromatin mark H3K9me3
243 (Figure 2a, Supplementary Figure 11); 4 'TN' sequence classes, which are enriched in transcription
244 elongation marks H3K36me3 or H3K79me2 (Figure 2a, Supplementary Figure 12); and finally, 7
245 'L' (low signal) sequence classes, which are not strongly enriched in any of the above marks
246 (Figure 2a). As a whole, the 40 sequence classes cover >97.4% of the genome (Supplementary
247 Figure 13).

248
249 Beyond classifying genomic sequences to sequence classes, we define sequence class scores to
250 provide a global and quantitative representation of sequence regulatory activities. This for the
251 first time allows us to (1) predict the regulatory activity for any sequence and (2) quantify the
252 changes in regulatory activity caused by any sequence variant. Sequence class scores summarize
253 predictions for all 21,907 chromatin profiles based on weights specific to each sequence class,
254 which are computed by projecting Sei predictions onto unit-length vectors that point to the center
255 of each sequence class. Sequences that score highly for a particular sequence class have high
256 predictions for the chromatin profiles associated with that class. Sequence class scores thus allow
257 for the quantification of the regulatory activity of any sequence, where the impact of a variant is
258 represented by the difference between the sequence class scores for the reference and alternative
259 alleles. Importantly, this capability is only allowed by modeling the sequence dependencies of
260 sequence class activities and cannot be directly obtained from chromatin profiling data alone.

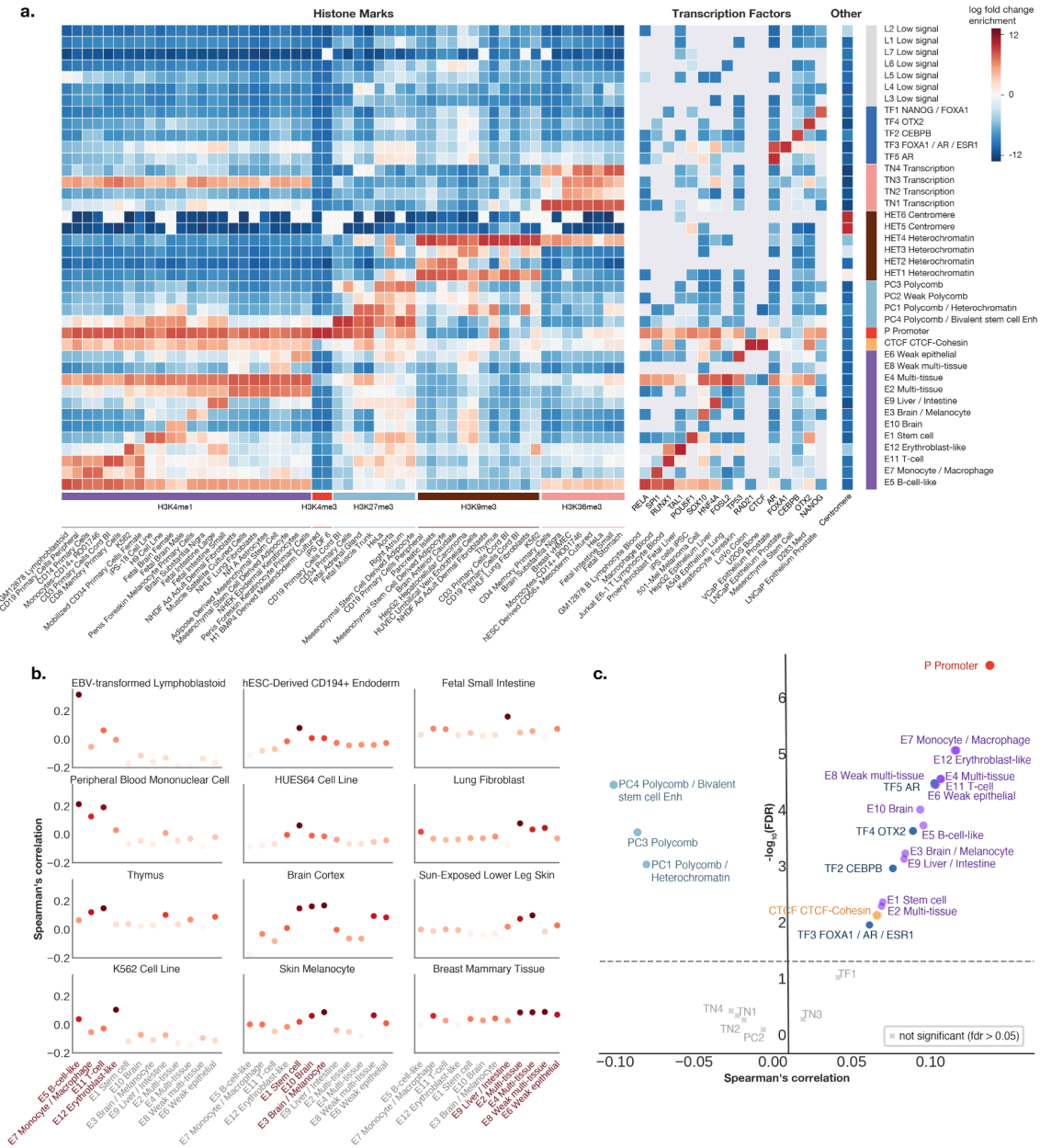
261

262 **Enhancer sequence classes predict tissue-specific gene expression**

263 The group of sequences that are likely most impactful to tissue-specific gene expression
264 regulation are the enhancer ('E') sequence classes, thus here we assessed the association of
265 enhancer sequence class scores with tissue-specific gene expression.
266

267 In the visualization of sequence regulatory activities, sequence classes with different cell type-
268 and tissue-specific enhancer activities are localized to distinct subregions (Figure 1b). 'E'
269 sequence classes capture both specific and broad enhancer activities. Based on enhancer mark
270 enrichment (Supplementary Figures 8, 9), E7 is specific for monocyte/macrophage, E11 is
271 specific for T-cell, E5 is specific for lymphoblastoid/B-cell-like cell lines, E9 is specific for liver
272 and intestine, E1 is specific for embryonic stem cells & induced pluripotent stem cells, and E10
273 and E3 are specific for brain (Figure 1, 2; all enrichments stated are significant with $p < 2.2e-16$,
274 Fisher's exact test, two-sided). In contrast, broad enhancer sequence classes can either
275 encompass enhancer activity in similar cell types across different tissues, such as fibroblast (E2)
276 and epithelial (E6) cell types (Supplementary Figures 8, 9), or encompass enhancer activity in
277 many different cell types; for example, E4 is enriched in fibroblast, muscle, astrocytes,
278 osteoblast, epithelial, and other cell types. Sequence class enhancer activities are also supported
279 by the enrichment of relevant chromatin states³ and DNase I hypersensitive sites²⁵ across tissues
280 and cell types (Supplementary Figures 14, 15). Consistent with their predicted enhancer
281 activities, the coverage of 'E' sequence class annotations within a 10kb window to transcription
282 start sites (TSS) are correlated with the differential expression patterns of these genes in the
283 corresponding cell types over the tissue-average (Figure 2b).
284

285 Since sequence class scores allow us to systematically predict the effects of variants on higher-
286 level regulatory functions, we can estimate whether a given variant diminishes, maintains, or
287 increases the enhancer activity of a sequence based on the difference between the sequence class
288 scores for the reference and alternative alleles. Evaluated on GTEx eQTL data²⁶, we found that
289 variants predicted to increase 'E' sequence class activity were significantly positively correlated
290 with higher gene expression, whereas those predicted to increase 'PC' sequence class activity were
291 significantly negatively correlated with gene expression--consistent with the expected repressive
292 role of 'PC' sequence class activities (Figure 2c). Moreover, when only analyzing fine-mapped
293 eQTLs²⁷ with high posterior inclusion probability (>0.95), we observed higher correlations with
294 overall comparable levels of significance (Supplementary Figure 16). Therefore, sequence classes
295 can distinguish the effects of variants on gene expression based on their consequences in regulatory
296 activities.



297
 298 **Figure 2. Sequence classes predict cell-type-specific regulatory activities and directional,**
 299 **expression-altering variant effects.** **a**, Sequence-class-specific enrichment of histone marks,
 300 transcription factors, and repeat annotations. Log fold change enrichment over genome-average
 301 background is shown in the heatmap. No overlap is indicated by the gray color in the heatmap. Top 1-2
 302 histone mark and TF annotation enrichments were selected for each sequence class. **b**, Enhancer sequence
 303 classes near transcription start sites are correlated with cell-type-specific gene expression in the applicable
 304 tissue or cell types (see Methods). The y-axis shows the Spearman correlation between the proportion of
 305 each sequence class annotation within 10kb of TSS and the tissue-specific differential gene expression
 306 (fold over tissue-average). **c**, Regulatory sequence-class-level variant effects are predictive of directional
 307 GTEx variant gene expression effects. The x-axis shows Spearman correlations between the predicted
 308 sequence-class-level variant effects and the signed GTEx variant effect sizes (slopes) for variants with

309 strong predicted effects near transcription start sites (Methods) and the y-axis shows the corresponding
310 log₁₀ p-values. All colored dots are above the Benjamini-Hochberg FDR < 0.05 threshold.

311

312

313 **Regulatory sequence classes are under evolutionary constraints**

314 Variants that alter regulatory activities of sequences often disrupt gene regulation and are
315 therefore expected to impact human health and disease. We tested this expectation by comparing
316 human population genome variant allele frequencies²⁸ based on the sequence class in which each
317 variant is located and the predicted variant effect on that sequence class. Indeed, we found that
318 variants localized in regulatory sequence classes (E-, P-, and CTCF-) have lower common
319 variant frequency than variants in other sequence classes, and therefore showed higher overall
320 negative selection constraint (Figure 3a, x-axis). More importantly, variants predicted to strongly
321 perturb regulatory sequence classes had significantly lower common variant frequencies than
322 variants that weakly perturb these classes (measured by bidirectional variant effect constraint,
323 Figure 3a y-axis, see also Figure 3b, Methods). This is therefore consistent with the hypothesis
324 that disruption of regulatory sequence class activities has a major negative impact on fitness,
325 which we refer to as a negative selection signature.

326

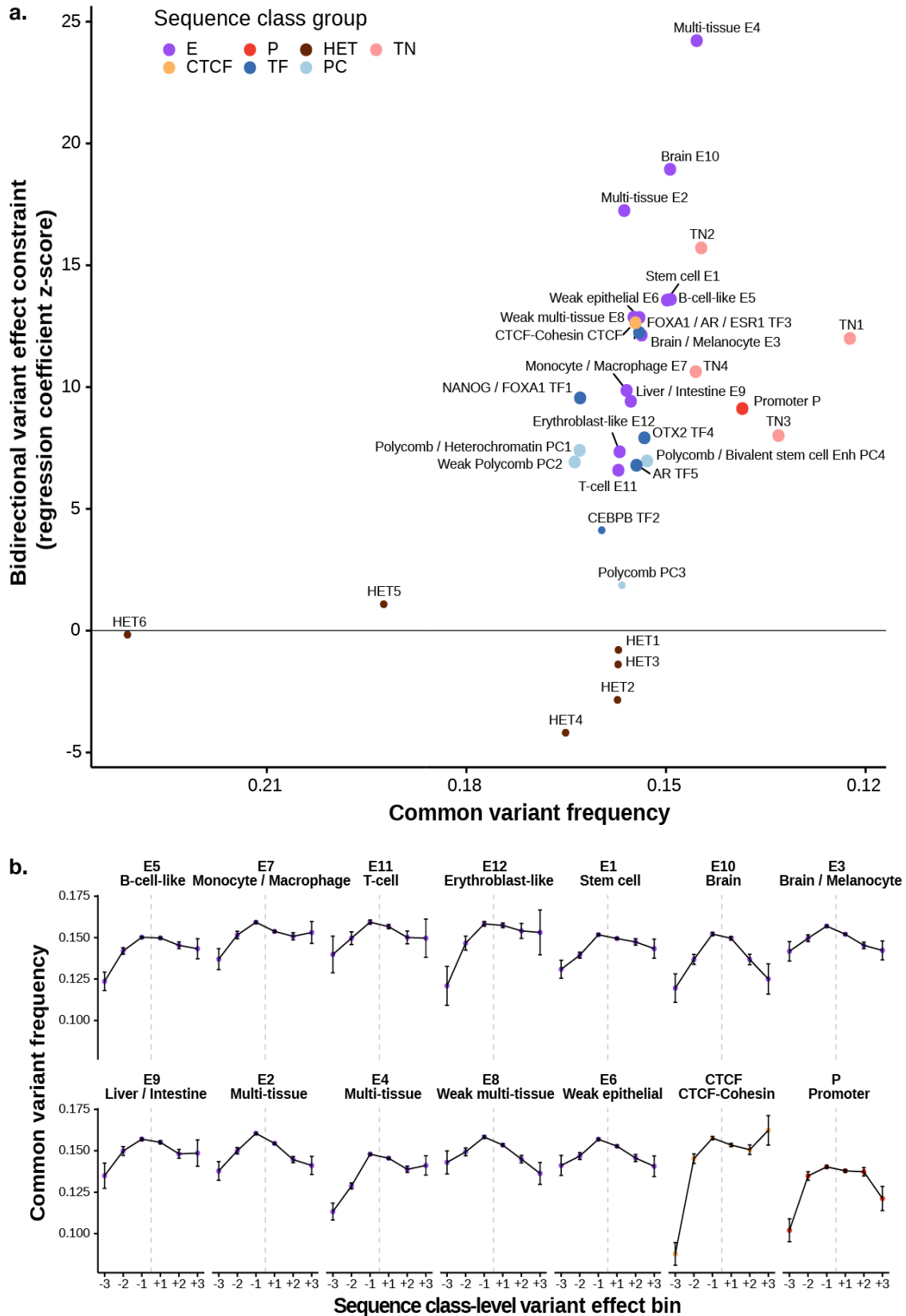
327 Specifically, we observed strong negative selection signatures for variants assigned to all E,
328 CTCF and P sequence classes (Figure 3). Multi-tissue enhancer sequence classes E4 and E2 and
329 the brain enhancer sequence class E10 show the strongest association of predicted sequence-
330 class-level variant effects and common variant frequencies. Notably, for the CTCF sequence
331 class, only negative variant effects--decreasing sequence class activity--appear to be under very
332 strong constraints, suggesting that CTCF sites are generally tolerant to positive effect mutations
333 that further increase CTCF binding. This is in contrast to the generally deleterious impact of both
334 increase and decrease of enhancer and promoter activities. As expected, TN sequence classes,
335 which overlap with protein-coding regions, are among the sequence classes with the lowest allele
336 frequency (Supplementary Figure 17).

337

338 In contrast, those assigned to HET, PC, TF, and L sequence classes generally did not show
339 strong negative selection signatures and had higher overall common variant frequencies
340 (Supplementary Figure 17). Importantly, this does not suggest that Polycomb or transcription
341 factors are inessential: the HET, PC, TF, and L classes generally do not show strong enhancer or
342 promoter histone mark enrichment in any cell type (with the exception of bivalent marks in stem
343 cells observed in PC4), and thus they are expected to play less major roles in gene expression
344 regulation. However, Polycomb-related regulation is likely critical for 'E' and 'P' sequence
345 classes, which are often Polycomb-repressed in some cell types but enhancers or promoters in
346 other cell types (Supplementary Figures 7-10). Similarly, we expect that TF binding plays a
347 central role in 'E' classes that are highly enriched in relevant TFs (Figure 2a, Supplementary File
348 3).

349

350 Therefore, sequence classes show distinct evolutionary constraints, and ‘E’ enhancer sequence
351 classes show the strongest bidirectional constraints. This suggests that both increases and
352 decreases of enhancer activity are expected to lead to deleterious effects on fitness, highlighting
353 the importance of precisely controlling gene expression.



355 **Figure 3. Variants with strong regulatory sequence class effects show negative selection signatures.**
356 **a**, Scatter plot for allele-frequency-based analysis of each sequence class. The x-axis shows 1 - common
357 variant frequency (allele frequency > 0.01) across all 1000 Genome variants per sequence class, and the
358 y-axis shows the bidirectional variant effect constraint z-score, which is computed based on logistic
359 regressions predicting common variant (allele frequency > 0.01) from sequence-class-level variant effect
360 score for both positive and negative effects (Methods). Sequence classes with significant (Bonferroni-
361 Hochberg FDR<0.05) bidirectional variant effect constraint are indicated with larger dots. ‘L’ sequence
362 classes are excluded due to lack of interpretation for their sequence-class-level variant effect scores. **b**,
363 Comparison of common variant frequencies for 1000 Genomes variants assigned to different sequence
364 classes and variant effect bins. The common variant threshold is >0.01 allele frequency across the 1000
365 Genomes population. Error bars show +/- 1 standard error (SE). The sequence-class-level variant effects
366 are assigned to 6 bins (+3: top 1% positive, +2: top 1%-10% positive, +1: top 10% -100% positive, -1: top
367 10% -100% negative, -2: top 1%-10% negative, -3: top 1% negative).

368

369

370 **Sequence classes elucidate the tissue-specific regulatory architecture of GWAS traits**

371 The population allele frequency analysis on sequence classes suggest that variants perturbing
372 regulatory sequence class activities are likely involved in human health and disease. Therefore,
373 to explore this hypothesis, we used GWAS data to delineate the genetic contribution of each
374 sequence class to diseases and traits.

375

376 Partitioned heritability from LD score regression (LDSR) has been a powerful tool for
377 understanding the genetic architecture of diseases and traits using GWAS summary statistics²⁹,
378 including identifying enrichment of disease heritability in regulatory elements^{29,30}. Previous
379 applications of LDSR use overlapping annotations,²⁹⁻³¹ which allows for the joint analysis of
380 heritability contribution across a wide range of annotations and has generated significant insight
381 into a wide range of GWAS studies; however, such analyses cannot unambiguously partition
382 heritability across annotations. Because sequence classes are both non-overlapping and cover
383 nearly the entire genome, they provide a clear and more easily interpretable picture of the
384 regulatory architecture of diseases and traits. To show this, we estimated the proportion of
385 heritability explained by each sequence class for 47 GWAS traits in UK Biobank (UKBB)^{18,32}
386 (Methods). Specifically, we applied LDSR and used a conservative estimate of the proportion of
387 heritability, subtracting one standard error and lower-bounding by 0. Our analysis of UKBB
388 GWAS revealed genetic signatures of sequence-class-specific regulatory functions (Figure 4,
389 Supplementary File 4).

390

391 Importantly, ‘E’ and ‘P’ sequence classes cover almost all classes that explain a high proportion
392 of heritability for GWAS traits and diseases--the same sequence classes inferred to be under
393 strong evolutionary constraints (Figure 3a, Supplementary File 4). We observed three main
394 groups of traits that share similar heritability composition signatures across sequence classes.
395 The first group is blood-related traits, which contains two subgroups of immune-related and non-
396 immune-related traits. The majority of heritability signals in blood-related traits are explained by

397 enhancer classes for the relevant cell type(s), such as monocyte/macrophage enhancer (E7) for
398 *Monocyte Count*, B-cell-like enhancer (E5) for *Auto Immune Traits*, and erythroblast-like (red
399 blood cell progenitor) enhancer (E12) for *Red Blood Cell Distribution Width*, which measures
400 the range of variation in red blood cell volume. Furthermore, autoimmune-related traits are
401 selectively associated with the immune cell type enhancer sequence classes E5 (B-cell like), E11
402 (T-cell), and E7 (monocyte/macrophage), while erythroblast-like enhancer E12 is specifically
403 linked to red-blood-cell-related traits. Therefore, sequence classes can dissect the cell-type-
404 specific regulatory architecture of traits and diseases with heritability decomposition, even
405 without relying on gene-level information.

406
407 Cognitive and mental traits (*Morning Person*, *Neuroticism*, *Smoking Status*, *Years of Education*,
408 *College Education*) have similar sequence-class-level heritability decompositions as well; for
409 this second group of traits, heritability was mostly explained by brain enhancer (E10 and E3) and
410 stem cell enhancer (E1) sequence classes. The link to E1 is consistent with our observation that
411 E1 was also moderately enriched for active enhancer mark H3K4me1 in brain cell types (Figure
412 2a, Supplementary Figure 7) and is positively correlated with gene expression in brain tissues
413 (Figure 2b).

414
415 The third group of traits is intriguingly diverse, including *Balding*, *Lung Forced Vital Capacity*,
416 *Waist-hip Ratio*, *Height*, and *Heel T-score*. The heritability of these traits are mostly explained
417 by multi-tissue enhancer classes (E4, E2, and E8), which show activity in epithelial cells,
418 fibroblast, muscle, and many other cell types. Enhancer activity across multiple tissues in the
419 body may explain the diverse phenotypes that are associated with these traits.

420
421 Beyond these three groups, there are a number of traits with unique heritability patterns that are
422 also linked to highly relevant sequence classes. For example, the *High Cholesterol* trait was most
423 associated with the liver and intestine enhancer sequence class (E9), which is consistent with the
424 physiology of cholesterol metabolism and known etiology of this condition³³. E9 was also linked
425 to red-blood-cell-related traits, in line with the role of liver in erythropoiesis.

426
427 Finally, the promoter sequence class P uniquely explained a sizable proportion of heritability in
428 nearly all traits, suggesting a near-universal involvement of promoter sequence variations in all
429 traits and diseases.



430
 431 **Figure 4. Sequence-class-based partitioning of GWAS heritability shows trait associations with**
 432 **tissue-specific regulation.** Partitioned genome-wide heritability in UKBB GWAS with all 40 sequence
 433 classes. The size of the dot indicates the proportion of heritability estimated from LDSR, which is
 434 conservatively estimated as one standard error below the estimated heritability proportion (bounded by 0).
 435 The color of the dot indicates the significance z-score of the fold enrichment of heritability relative to the
 436 proportion of all SNPs assigned to the sequence class (bounded by 0). Colored boxes indicate traits
 437 associated with blood (red), brain (green), multiple tissues (blue) and promoters (orange).

438
 439
 440 We next assessed whether our new sequence classes could explain GWAS heritability beyond
 441 that explained by annotations discovered in prior studies. To this end, we performed LDSR
 442 analysis with our whole genome annotations of sequence classes conditioned on an up-to-date set
 443 of previously identified baseline annotations (v2.2,
 444 <https://alkesgroup.broadinstitute.org/LDSCORE/>). We uncovered 83 significant sequence-class-

445 trait associations with a corrected p-value cutoff of <0.05 (Supplementary File 5). 70% of all
446 UKBB GWAS traits and 9/13 of the E and P sequence classes have at least one significant
447 association after multiple hypothesis testing correction (Supplementary File 5). This finding
448 suggests that sequence classes can identify extensive new regulatory signals that enrich GWAS
449 interpretation.

450

451 **Disease mutations are predicted to disrupt the activities of sequence classes**

452 Sequence-class-level effects enable the prediction of specific regulatory mechanisms at the
453 individual, pathogenic mutation level. To showcase our framework's capability to predict the
454 mechanisms of individual mutations, we used Sei to predict the direction and magnitude of
455 sequence-class-level mutation effects for all 853 regulatory disease mutations from the Human
456 Gene Mutation Database (HGMD)³⁴. For systematic classification and quantification of these
457 mutations, we assign each mutation to an affected sequence class based on its mutation effects
458 (the sequence class with the strongest score change) and the sequence that it alters (Methods).

459

460 Overall, the average variant effect score of disease mutations is 4.2x larger than the de novo
461 mutations in healthy individuals (0.903 vs 0.217, $p < 2.2e-16$, Wilcoxon rank-sum test two-sided,
462 max absolute effect across sequence classes) and 6.5x larger than the 1000 Genomes common
463 variants with $AF > 0.01$ (0.903 vs 0.139, $p < 2.2e-16$). Here we focus on analyzing the mutations
464 with the strongest predicted effects (>1.1 , $n=138/853$), where predicted effect refers to the
465 variant effect of the assigned sequence class for each mutation (Figure 5, Supplementary Figure
466 18). Because sequence-class-level variant effects are directional--that is, predicting whether the
467 alternative allele increases or decreases sequence-class level activity--we are able to discover that
468 while the majority (~80%) of pathogenic mutations with strong predicted effects are predicted to
469 decrease sequence class activity, the remaining 20% of HGMD pathogenic mutations are
470 predicted to increase sequence class activity. Moreover, perturbations to E-, P-, and CTCF-
471 classes make up >99% of the mutations with strong predicted effects on sequence class activity
472 (Supplementary File 6): 44.9% are predicted to affect tissue-specific E sequence classes, 38.4%
473 are predicted to affect the P promoter sequence class, and interestingly, 15.9% are predicted to
474 affect the CTCF-cohesin sequence class (Methods).

475

476 We found that almost all mutations with strong predicted effects in cell-type-specific E sequence
477 classes contributed to diseases relevant to that same cell type (Figure 5, Supplementary File 6)--
478 for most of these mutations, the nearby gene is known to be relevant to the disease but the
479 molecular mechanisms of regulatory disruption is unknown. For example, mutations causing
480 Protein C deficiency and Hemophilia B, two diseases characterized by the deficiency of specific
481 plasma proteins produced in the liver (protein C and coagulation factor IX, respectively), are
482 predicted to decrease E9 liver/intestine sequence class activities. Blood cell-type-specific
483 enhancer sequence classes are disrupted in distinct blood-related diseases and deficiencies
484 relevant to the corresponding cell type: the E12 erythroblast-like enhancer sequence class is

485 disrupted in red blood cell-specific diseases such as pyruvate kinase deficiency, erythropoietic
486 porphyria, delta-thalassemia, and beta-thalassemia; the E7 monocyte/macrophage-like sequence
487 class is disrupted in monocyte and macrophage-related chronic granulomatous disease; and the
488 E5 B-cell-like enhancer sequence class is disrupted in X-linked agammaglobulinemia, a
489 functional deficiency of B-cell. For developmental diseases, such as preaxial polydactyly
490 triphalangeal thumb and radial ray deficiency and triphalangeal thumb-polysyndactyly
491 syndrome, the E1 embryonic stem cell-specific enhancer sequence class is predicted to be
492 disrupted by mutations in a known distal enhancer of Sonic Hedgehog (SHH) (chr7:156583951
493 G>A³⁵, chr7:156583949 G>C³⁶), a gene that plays a crucial role in the positioning and growth of
494 limbs, fingers, and toes during development.

495
496 In addition, 38% of the regulatory mutations with strong predicted effects affect the activity of
497 the promoter sequence class P, including a hypercholesterolemia mutation near the LDLR gene
498 (chr19:11200089 C>T³⁷), a microcephaly & developmental delay mutation near the PIGY gene
499 (chr4:89444948 C>T³⁸), and a retinoblastoma mutation near the RB1 gene (chr13:48877851
500 G>T³⁹). The high proportion of mutations perturbing the P sequence class likely reflects both the
501 critical role of promoters in diseases and the emphasis on promoter-proximal mutations in past
502 studies.

503
504 While the mutations we've discussed thus far are negative effect mutations which decrease
505 sequence class activity, 20% of HGMD pathogenic mutations are predicted to increase sequence
506 class activity. Indeed, these mutations included many known gain-of-function mutations, which
507 validated our predictions. The highest increase in sequence class activity was observed for a
508 mutation (chrX:73072592 G>C) near the XIST gene that skews X-inactivation of the mutant
509 chromosome in females⁴⁰; this mutation was predicted to increase the activity of the CTCF
510 sequence class and has been experimentally validated to increase CTCF binding⁴¹. Similarly,
511 positive effect predictions for 'E' and 'P' sequence classes were also validated by previously
512 studied mutations: an alpha-thalassemia mutation near the HBM gene (chr16:209709 T>C⁴²)
513 known to create a GATA1 binding site and increase intergenic transcription was predicted to
514 increase the activity of the erythroblast-specific E12 sequence class, and a TERT gene mutation
515 found in individuals with familial melanoma (chr5:1295161 T>G⁴³) was predicted to increase the
516 activity of P. Beyond this, many mutations predicted to have strong positive effects were not
517 previously understood. For example, a mutation near the HBG1 gene (chr11:5271262 A>G⁴⁴)
518 that causes persistence of fetal hemoglobin is also predicted to increase the activity of the
519 erythroblast-specific E12 sequence class. Previously, this mutation was known to create an
520 ATGCAAAT octamer⁴⁴ that matches the POU family transcription factor motif, but its
521 functional consequences were unclear.

522
523 Notably, even though pathogenic mutations from prior genetics studies are subjected to selection
524 bias, the observation of pathogenic mutations with strong impacts on E-, P-, and CTCF-

540 **Figure 5. Disease regulatory mutations are predicted to disrupt promoter, CTCF, and tissue-**
541 **specific enhancer sequence classes.** Sequence-class-level mutation effects of pathogenic noncoding
542 HGMD mutations are plotted. A polar coordinate system is used, where the radial coordinate indicates the
543 sequence-class-level effects. Each dot represents a mutation, and mutations inside the circle are predicted
544 to have positive effects (increased activity of sequence class), while mutations outside of the circle are
545 predicted to have negative effects (decreased activity of sequence class). Dot size indicates the absolute
546 value of the effect. Mutations are assigned to sequence classes based on their sequences and predicted
547 effects (Methods). Within each sequence class, mutations are ordered by chromosomal coordinates. The
548 associated disease and gene name are annotated for each mutation, and only the strongest mutation is
549 annotated if there are multiple mutations associated with the same disease, gene, and sequence class.

550

551

552 **Discussion**

553 We developed a genome-wide sequence-based map of regulatory activities using sequence
554 classes, a vocabulary for genomic sequence activities discovered using a data-driven, systematic
555 method. Our deep-learning-based framework uses a compendium covering 21,907 publicly
556 available cis-regulatory profiles and the whole genome sequence to create a mapping from any
557 sequence to a comprehensive set of sequence classes. This provides a global sequence-based
558 view of sequence regulatory activities and allows for the quantitative prediction of variant effects
559 on sequence class activities. Sequence classes are a concise vocabulary of regulatory activities
560 that is interpretable, quantifiable, and easily analyzed globally (across all sequence classes) and
561 individually. To our knowledge, it is the first such attempt to systematically map regulatory
562 activities from any sequence.

563

564 We demonstrated that E- and P- sequence classes are strongly enriched in trait and disease
565 GWAS heritability and under evolutionary constraints. Importantly, sequence classes provide
566 insights into the mechanisms of individual pathogenic mutations by predicting effects on the
567 function of tissue-specific enhancers, promoter activity, and long-range genome interactions (e.g.
568 CTCF-cohesin sequence class). Using sequence-class-level variant effect predictions, we linked
569 many pathogenic mutations to tissue-specific regulatory changes in the relevant tissues. These
570 predictions point to potential mechanisms that can be experimentally tested in the future.

571

572 Sequence classes leverage a sequence model trained on most publicly available cis-regulatory
573 profile data; however, there remains substantial space for improvement as more data becomes
574 available. For example, we are still lacking data for many cell types, developmental stages,
575 transcription factors, and combinations of chromatin targets measured in new cell types or
576 conditions. More data that covers currently undercharacterized cell types and developmental
577 stages will likely enable the identification of still more cell-type-specific and developmental
578 stage-specific sequence classes, defining sequence classes with increasingly fine-grained
579 regulatory resolution. Furthermore, development of new computational methods to define, for
580 example, hierarchical or combinatorial representations of sequence classes may be needed to

581 make such fine-grained classes easy to interpret and use. Because interpretability and robustness
582 were our major goals in designing sequence classes, we chose to use clustering to generate the
583 sequence classes and a linear projection step to compute corresponding scores. It is conceivable
584 that a more expressive model such as an end-to-end neural network can further improve
585 sequence class predictions, and we expect that increasing the expressiveness of the model while
586 maintaining interpretability and robustness will be an interesting future challenge.

587
588 This work demonstrates the potential of sequence classes to discover regulatory disruptions in
589 human diseases, through both the aggregation of genome-wide variant association signals and
590 prediction of the impact of individual mutations. We provide sequence classes and the Sei model
591 as a resource for further research into understanding the regulatory genetic landscape of human
592 health and diseases. Our framework is applicable to any variant, regardless of whether it is
593 common, rare, or never previously observed, and we expect it to be a powerful tool for
594 understanding the mechanistic effects of noncoding mutations in human health.

595

596 **Methods**

597

598 **Training data**

599 21,907 cis-regulatory profiles in peak format were compiled from the processed files of the
600 Cistrome⁵, ENCODE³, and Roadmap Epigenomics projects⁴. The Cistrome Project, which
601 systematically processed publicly available cis-regulatory profiles, contributed the majority of
602 the profiles predicted in Sei (19,905). We excluded profiles from Cistrome with less than 1000
603 peaks. Genome sequences are from the GRCh38/hg38 human reference genome. The full list of
604 cis-regulatory profiles is available in Supplementary File 1.

605

606 **Deep learning sequence model training**

607 The Sei model is trained to predict 21,907 transcription factor binding, histone marks, and DNA
608 accessibility from cis-regulatory profile peaks at the center of 4kb sequences.

609

610 The model architecture is composed of three sequential sections: 1) a convolutional network with
611 dual linear and nonlinear paths, 2) residual dilated convolution layers, 3) spatial basis function
612 transformation and output layers. A detailed specification of the model is available in
613 Supplementary File 7 and in the code repository (<https://github.com/FunctionLab/sei-framework>,
614 downloadable from <https://doi.org/10.5281/zenodo.4906996>). In the convolutional architecture,
615 we introduced a new design composed of both linear and nonlinear convolution blocks. The
616 nonlinear blocks are composed of convolution layers and rectified linear activation functions
617 (ReLU), similar to regular convolutional networks. The linear blocks have the same structure as
618 the nonlinear blocks but do not include activation functions to facilitate learning of linear
619 dependencies. Each nonlinear block is stacked on top of a linear block with a residual connection
620 adding the input of the nonlinear block to the output, allowing the computation to go through

621 either the linear or nonlinear path. Dilated convolutional layers with residual connections further
622 expands the receptive fields without reducing spatial resolution. Finally, spatial basis functions
623 are used to reduce dimensionality of the spatial dimension while preserving the capability to
624 discriminate spatial patterns of sequence representations. Specifically, in the Sei model, a B-
625 spline basis matrix (256x16) with 16 degrees-of-freedom across 256 uniformly-spaced spatial
626 bins is generated and multiplied with the convolutional layers output to reduce the 256 spatial
627 dimensions to 16 spline basis function dimensions. After the spline basis function
628 transformation, a fully-connected layer and an output layer are used for integrating information
629 across the whole sequence and generating the final 21,907-dimensional predictions.

630
631 Training, validation, and testing datasets are specified by different sets of chromosomes in the
632 hg38 genome (holding out chromosome 8 and 9 for the test set and chromosome 10 for the
633 validation set), and samples drawn uniformly across the hg38 genome for these partitions,
634 excluding regions specified in the ENCODE blacklist⁴⁵. For training, we sampled training
635 sequences and their labels on-the-fly from the training set of chromosomes using Selene⁴⁶. As a
636 result, almost all training samples are drawn from unique genomic intervals with distinct start
637 and end positions to reduce overfitting during the training process. For each 4kb region, a
638 21,907-dimensional binary label vector is created for the 21,907 cis-regulatory profiles based on
639 whether the center basepair overlaps with a peak in each of the profiles. The model is
640 implemented in PyTorch and trained with Selene. A detailed training configuration file is
641 available at <https://github.com/FunctionLab/sei-framework/blob/main/train/train.yml>.

642

643 **Model performance**

644 We computed the AUROC and AUPRC for all cis-regulatory profiles predicted by Sei on the test
645 holdout dataset, excluding profiles that had fewer than 25 positive samples in the test set.
646 Additionally, to assess the correlation structure of the predictions, we compared the rank-
647 transformed pairwise Spearman's rank correlations for the predicted cis-regulatory profiles to the
648 pairwise correlations for the true labels (peak calls provided in Cistrome DB).

649

650 The model performance comparison between DeepSEA and Sei is computed on the 2,002 cis-
651 regulatory profiles from Roadmap and ENCODE that both DeepSEA and Sei predict. Because
652 both models have the same chromosomal test holdout (chr8 and chr9), we use the regions
653 specified in the DeepSEA test holdout set to create a common test dataset of sequences and
654 labels on which to evaluate the models.

655

656 **Sequence classes**

657 We selected 30 million genomic positions that uniformly tile the genome with 100bp step size
658 and then computed Sei predictions for 4kb sequences centered at each position. Sequences
659 overlapping with ENCODE blacklist regions⁴⁵ or assembly gaps ("N"s) are removed. To process
660 the 30 million x 21,907 predictions matrix, the dimensionality is first reduced with principal

661 component analysis (PCA). The PCA transformations were fitted with incremental PCA using a
662 batch size of 1,000,000 for one pass of the whole dataset, and genomic positions were randomly
663 assigned to batches. The top 180 principal components, scaled to unit variance, were used for
664 constructing a nearest neighbor graph where each node is connected to its k-nearest neighbors by
665 Euclidean distance (k=14). Louvain community clustering with default parameters was applied
666 to the nearest neighbor graph with the python-louvain package, which resulted in 61 clusters. We
667 refer to the largest 40 clusters as sequence classes and exclude the remaining (smallest) 21
668 clusters, which constitute <2.6% of the genome, from our analyses due to their size. These 21
669 clusters mainly display Low signal or Heterochromatin like enrichment (Supplementary Figure
670 19). We refer to this cluster assignment to sequence classes at 100bp resolution as sequence class
671 annotations. We visualized the genome-wide predictions by computing UMAP embedding with a
672 subsample of PCA-transformed Sei predictions of 30 million sequences, and then fine-tuned the
673 visualization with OpenTSNE. The detailed procedures are available in our code repository
674 (<https://github.com/FunctionLab/sei-manuscript>).

675

676 **Sequence class scores**

677 Each sequence class is represented as a unit vector in the 21,907-dimensional cis-regulatory
678 profile space, in the direction of the average prediction of all sequences assigned to this sequence
679 class among the 30 million. In more formal notation, the vector for sequence class i is $v_i =$

680 $\frac{\overline{p_{s \in \text{Sequence class } i}}}{\|p_{s \in \text{Sequence class } i}\|_2}$, where p_s represents the 21,907-dimensional Sei prediction for sequence s .

681 Each Sei prediction can then be projected onto any sequence class vector to obtain a sequence
682 class-level representation of the prediction, which we call sequence class score or $score_{s,i} = p_s \cdot$
683 v_i^T . In addition, predicted sequence-class-level variant effects are represented by the difference
684 between the sequence class scores of the sequences carrying the reference allele and the
685 alternative allele, or $score_{v,i} = score_{alt,i} - score_{ref,i}$. To better represent predicted variant
686 effects on histone marks, it is necessary to normalize for the nucleosome occupancy (e.g. loss-of-
687 function mutation near TSS can decrease H3K4me3 modification level while increasing
688 nucleosome occupancy, resulting in an overall increase in observed H3K4me3 quantity).
689 Therefore, for variant effect computation, we use the sum of all histone profile predictions as an
690 approximation to nucleosome occupancy and adjust all histone mark predictions to remove the
691 impact of nucleosome occupancy change (non-histone mark predictions are unchanged):

692
$$p^{hm*} = p^{hm}_{ref} \frac{\sum_k p^{hm^k}_{ref} + \sum_k p^{hm^k}_{alt}}{\sum_k p^{hm^k}_{ref}}; p^{hm*} = p^{hm}_{alt} \frac{\sum_k p^{hm^k}_{ref} + \sum_k p^{hm^k}_{alt}}{\sum_k p^{hm^k}_{alt}}$$

693 where $\sum_k p^{hm^k}_{ref}$ represents the sum over all histone mark predictions (among 21907-
694 dimensions of a prediction) for the reference allele. We generally exclude Low Signal sequence
695 classes in sequence-class-level variant effect analyses because they lack an intuitive biological
696 interpretation.

697

698 **Sequence class enrichment of chromatin profiles and genome annotations**

699 We computed the log fold change enrichment of various chromatin profiles and genome
700 annotations for each sequence class based on sequence class annotations (described above, see
701 ‘Sequence classes’). Log fold change enrichment is computed by taking the log ratio of the
702 proportion of a sequence class intersecting with the annotation versus the background proportion
703 of the annotation, where we consider all regions assigned to any sequence class. We computed
704 enrichment for all 21,907 profiles predicted by Sei, filtered the chromatin profiles for each
705 sequence class to only those having Benjamini-Hochberg corrected p-values (Fisher’s exact test,
706 two-sided) below $2.2e-16$, and selected the top 25 profiles based on log fold change enrichment.
707 Cistrome Project profile enrichment is computed over 2 million random genomic positions.

708
709 The annotation of centromere repeats is obtained from the UCSC RepeatMasker track, and
710 annotations of histone marks over multiple cell types are obtained from the Roadmap
711 Epigenomics project--enrichments for both of these sets of annotations are computed over the
712 entire genome. In addition, we obtained ChromHMM chromatin states from ENCODE³ and
713 tissue and cell-type-specific DHS vocabulary from²⁵.

714

715 **Enhancer sequence class correlations with cell-type-specific gene expression**

716 Tissue expression profiles are from GTEx²⁶, Roadmap Epigenomics⁴, and ENCODE³ and
717 transformed to log-scale RPKM (reads per kilobase per million reads mapped) scores as
718 previously described¹³ and normalized by tissue-average. Specifically, a pseudocount was added
719 before log transformation (0.0001 for GTEx tissues, which are averaged across individuals, and
720 0.01 for Roadmap and ENCODE tissues). After log transformation, the average scores across
721 tissues were subtracted for each gene; as a result, the processed scores represent log fold change
722 relative to tissue-average.

723

724 Gene-wide expression prediction is evaluated on sequence class annotations (from Louvain
725 community clustering) for positions within +/-10kb of the TSSs for these genes. For each
726 enhancer sequence class and tissue, we compute the Spearman correlation between the sequence
727 class annotation coverage and gene expression.

728

729 **Correlation between regulatory sequence class variant effects and directional eQTL 730 variant effect sizes**

731 We collected the eQTLs within +/-5kb of gene TSSs from GTEx v8, combined across all GTEx
732 tissues, and computed the Spearman correlation between the top 15k variant effect predictions
733 for each sequence class and the eQTL variant effect sizes (averaged across multiple tissues if the
734 variant is an eQTL in multiple tissues). The p-values are derived from the Spearman’s rank
735 correlation test (two-sided) and BH correction is applied. Low Signal and Heterochromatin
736 sequence classes are excluded from this analysis due to lack of interpretation for their variant
737 effect scores in this context.

738

739 Additionally, we collected fine-mapped GTEx eQTLs from eQTL Catalogue²⁷ and obtained
740 sequence class scores for eQTLs with posterior inclusion probability > 0.95. Variants are
741 assigned to sequence classes based on the sequence class annotation for the reference genome
742 (i.e. variants are not further selected based on variant effect predictions). For each sequence
743 class, we computed the Spearman correlation between the sequence class scores and the eQTL
744 variant effect sizes in the same way we describe above.

745

746 **Evolutionary constraints on variant effects**

747 We computed sequence-class-level variant effects for all 1000 Genomes project phase 3
748 variants²⁸. Variants are assigned to sequence classes based on the 100bp resolution genome-wide
749 assignment derived from Louvain community clustering as described above. For each sequence
750 class we divide variants into 6 bins based on their effects in the same sequence class as
751 illustrated in Figure 3, and summarize common variant (AF>0.01) frequencies in each bin by
752 mean and standard error of the mean. We also estimated statistical significance of allele
753 frequency dependency on sequence-class-level variant effects. For each sequence class, we
754 applied logistic regression separately for positive effect and negative effect variants, to predict
755 common variants (AF>0.01) from the absolute value of sequence-class-level variant effect score,
756 and obtained the significance z-score of the regression coefficient of variant effect. The
757 bidirectional evolutionary constraint z-score is defined as the negative value of the combined z-
758 scores from positive and negative effect variants with Stouffer's method.

759

760 **Partitioning GWAS heritability by sequence classes**

761 UKBB GWAS summary statistics were obtained from¹⁸. To study the association of sequence
762 class genome annotation and sequence class variant effects and trait heritability, we performed
763 partitioned heritability LD score regression (LDSR) as described in²⁹. To partition the
764 heritability as sums of heritability explained by each sequence class, we run LDSR with only
765 sequence class annotations and a baseline all-ones annotation. We obtained the estimated
766 proportion of h^2 explained by each sequence class and its standard error with LDSR as
767 implemented in <https://github.com/bulik/ldsc>. As the estimated proportions can have high
768 variance or even be negative (the true value of heritability explained can only be non-negative),
769 we use a robust and conservative estimator which is the estimated proportion of h^2 subtracted by
770 one standard error, then lower-bounded by zero (the standard error of the estimated proportion of
771 h^2 explained is given by LDSR and estimated with the block jackknife procedure as described in
772 ²⁹).

773

774 To assess the contribution of sequence classes to explaining additional heritability when
775 conditioned on known baseline annotations, we also run LDSR with the baseline annotations
776 (v2.2, <https://alkesgroup.broadinstitute.org/LDSCORE/>). The p-values are derived from the
777 coefficient z-score, and BH correction is applied.

778

779 **Sequence class-level variant effect analysis of noncoding pathogenic mutations**

780 We obtained all mutations assigned “DM” and “regulatory” annotation in the Human Gene
781 Mutation Database (HGMD) database (2019.1 release). RMRP gene mutations are excluded
782 because they are likely pathogenic due to impacting RNA function instead of regulatory
783 perturbations, despite being annotated to the regulatory category in HGMD. For every mutation,
784 we predicted the sequence class scores for both the reference and the alternative allele and
785 computed the sequence-class-level variant effect as the predicted scores for the alternative allele
786 subtracting the scores for the reference allele. To provide an overview of sequence-class level
787 effects of human noncoding pathogenic mutations, mutations are first assigned to sequence
788 classes based on the sequence class annotations of the mutation position. For mutations with a
789 strong effect in a different sequence class than the originally assigned sequence class (absolute
790 value higher than the original sequence class by >1 absolute difference and >2.5 fold relative
791 difference), we reassign the mutation to the sequence class with the strongest effects.

792

793 **Code and data availability**

794 The Sei framework code is provided in <https://github.com/FunctionLab/sei-framework>, and the
795 model and associated data files downloadable by following the instructions in the GitHub
796 repository. Code and data for the manuscript results are available at
797 <https://github.com/FunctionLab/sei-manuscript>.

798

799 **Acknowledgements**

800 The authors acknowledge all members of the Troyanskaya lab for helpful discussions. This work
801 was performed using the high-performance computing resources, supported by the Scientific
802 Computing Core, at the Flatiron Institute and the Terascale Infrastructure for Groundbreaking
803 Research in Science and Engineering high-performance computer center at Princeton University.
804 K.M.C. is supported by the National Science Foundation Graduate Research Fellowship Program
805 (NSF-GRFP). O.G.T. is supported by National Institutes of Health grant nos. R01HG005998,
806 U54HL117798 and R01GM071966, U.S. Department of Health and Human Services grant no.
807 HHSN272201000054C and Simons Foundation grant no. 395506. O.G.T. is a senior fellow of
808 the Genetic Networks program of the Canadian Institute for Advanced Research. J.Z. is
809 supported by the Cancer Prevention and Research Institute of Texas grant (RR190071), National
810 Institutes of Health grant DP2GM146336, and the UT Southwestern Endowed Scholars program.

811

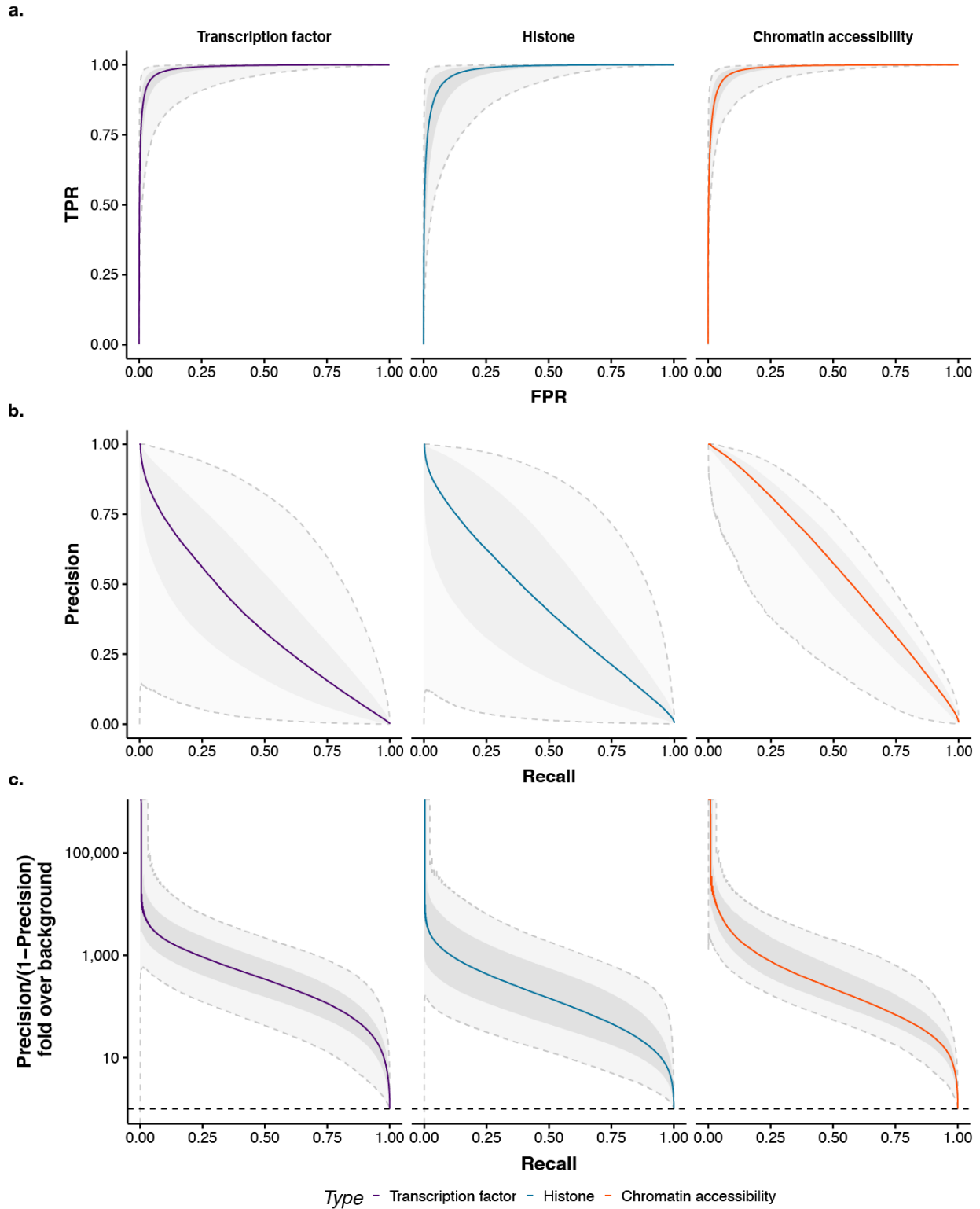
812 **Author Contributions**

813 K.M.C. and J.Z. conceived the Sei framework, developed the computational methods, and
814 performed the analyses. A.K.W. developed the Sei web server. K.M.C., J.Z., and O.G.T. wrote
815 the manuscript.

816

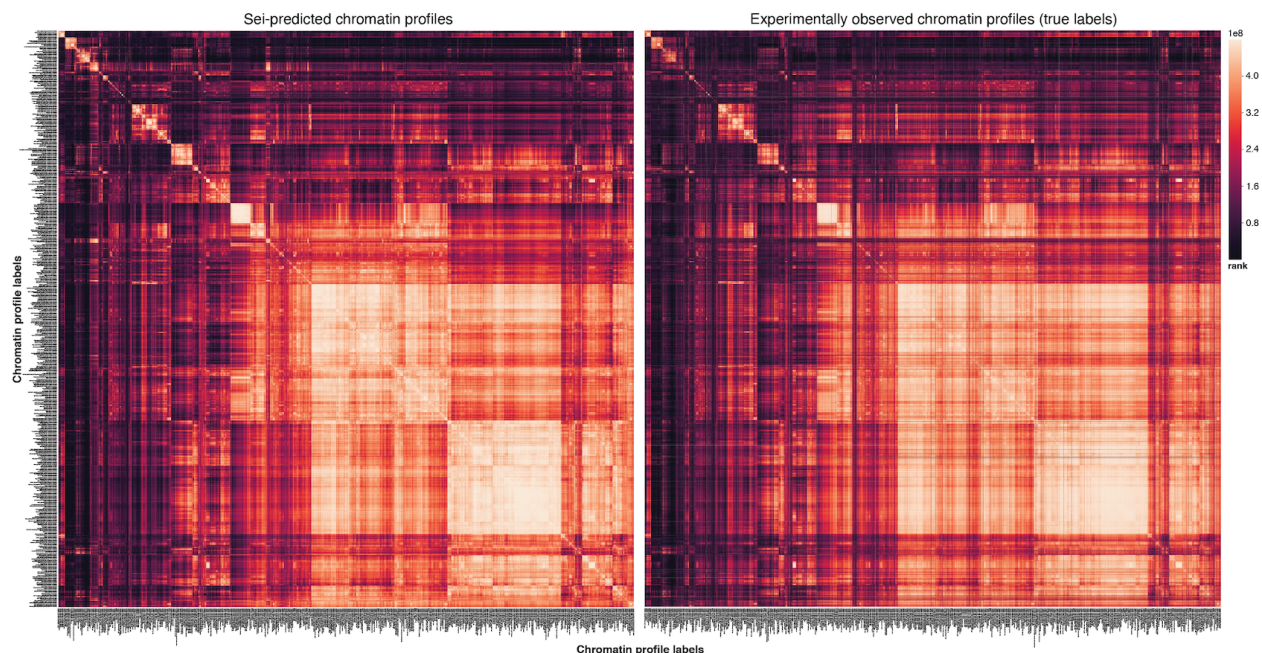
817

818

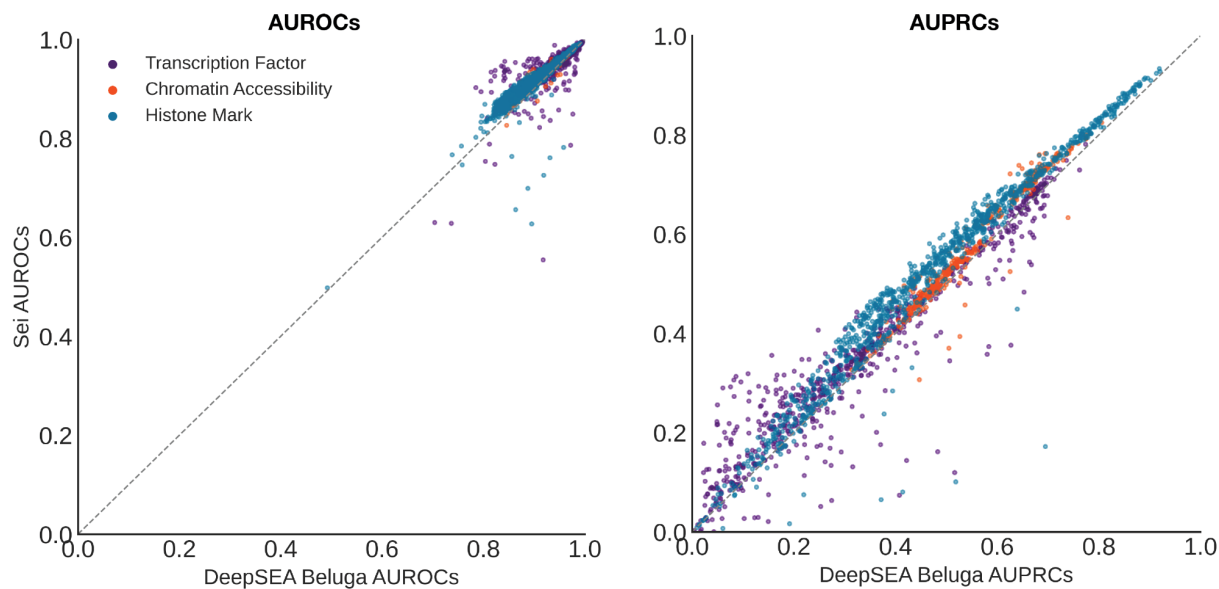


824
825
826
827

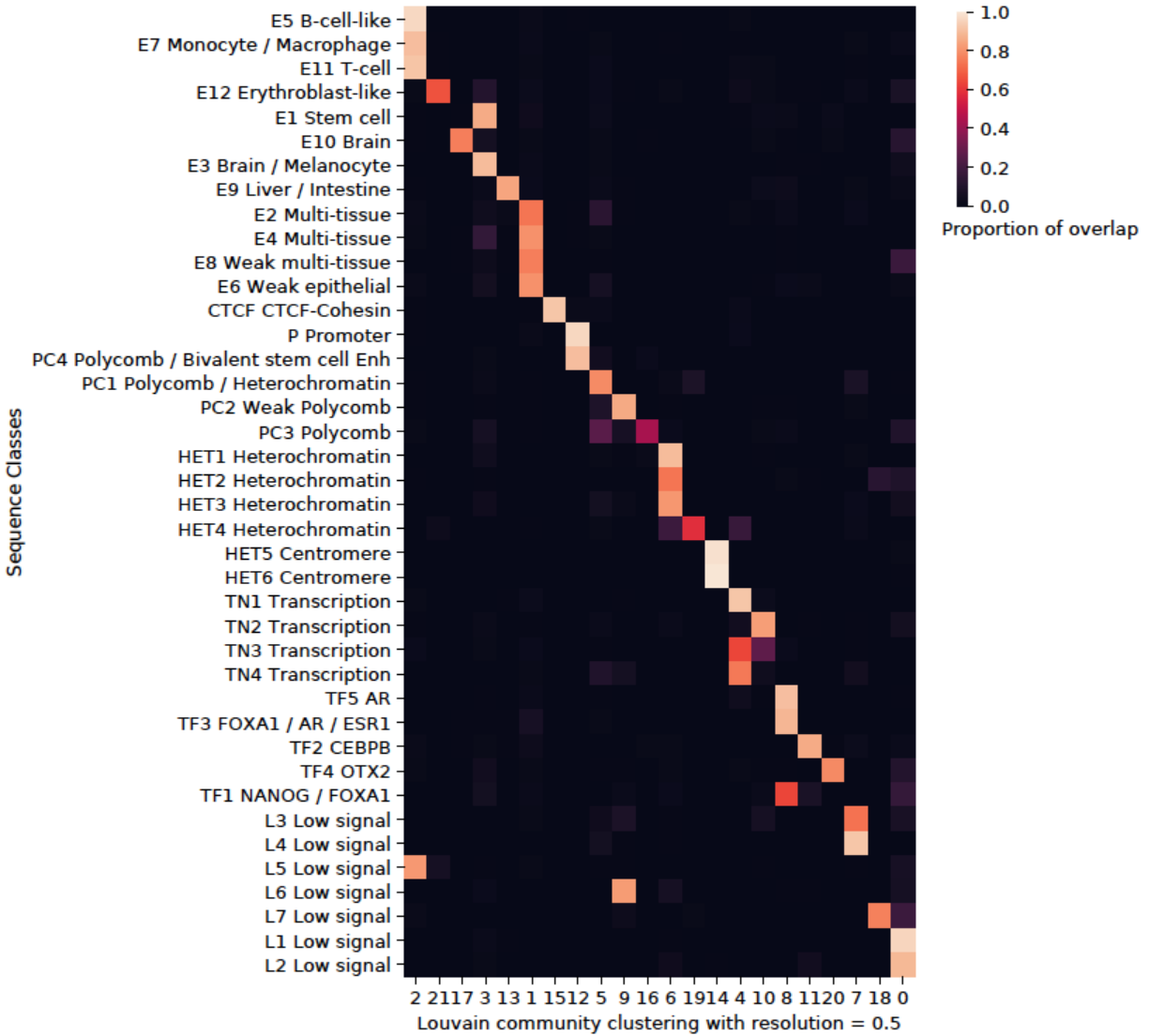
Supplementary Figure 2. Sei model performance on predicting 21907 cis-regulatory profiles on holdout chromosomes.



828
829 **Supplementary Figure 3. Visualizing the rank-transform of pairwise Spearman correlations for the**
830 **21,907 cis-regulatory profiles in Sei.** Sei model predictions share a highly similar correlation structure
831 with the experimental observations.
832



833
834 **Supplementary Figure 4. Sei model performance comparison with DeepSEA.** Performance on the
835 shared 2002 DeepSEA “Beluga” (2018) cis-regulatory profiles are compared.
836



837

838

Supplementary Figure 5. Comparison of sequence classes and Louvain community clustering with

839

resolution = 0.5. For each sequence class, the proportion overlap was computed between sequence

840

classes and a lower resolution clustering for Louvain community clustering. The lower resolution

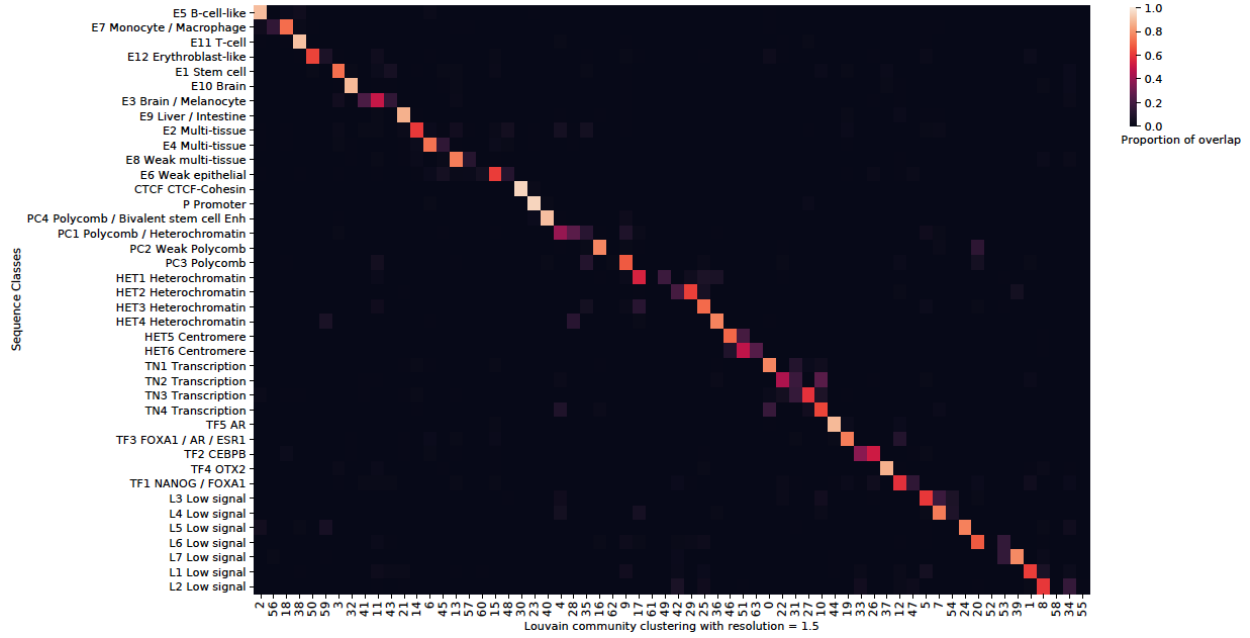
841

clustering is largely consistent with the original sequence classes, with some clusters combining several

842

related enhancer sequence classes into one.

843



844
845 **Supplementary Figure 6. Comparison of sequence classes and Louvain community clustering with**
846 **resolution = 1.5.** For each sequence class, the proportion overlap was computed between sequence
847 classes and a higher resolution clustering for Louvain community clustering. The higher resolution
848 clustering closely resembles the current sequence class clusters.
849



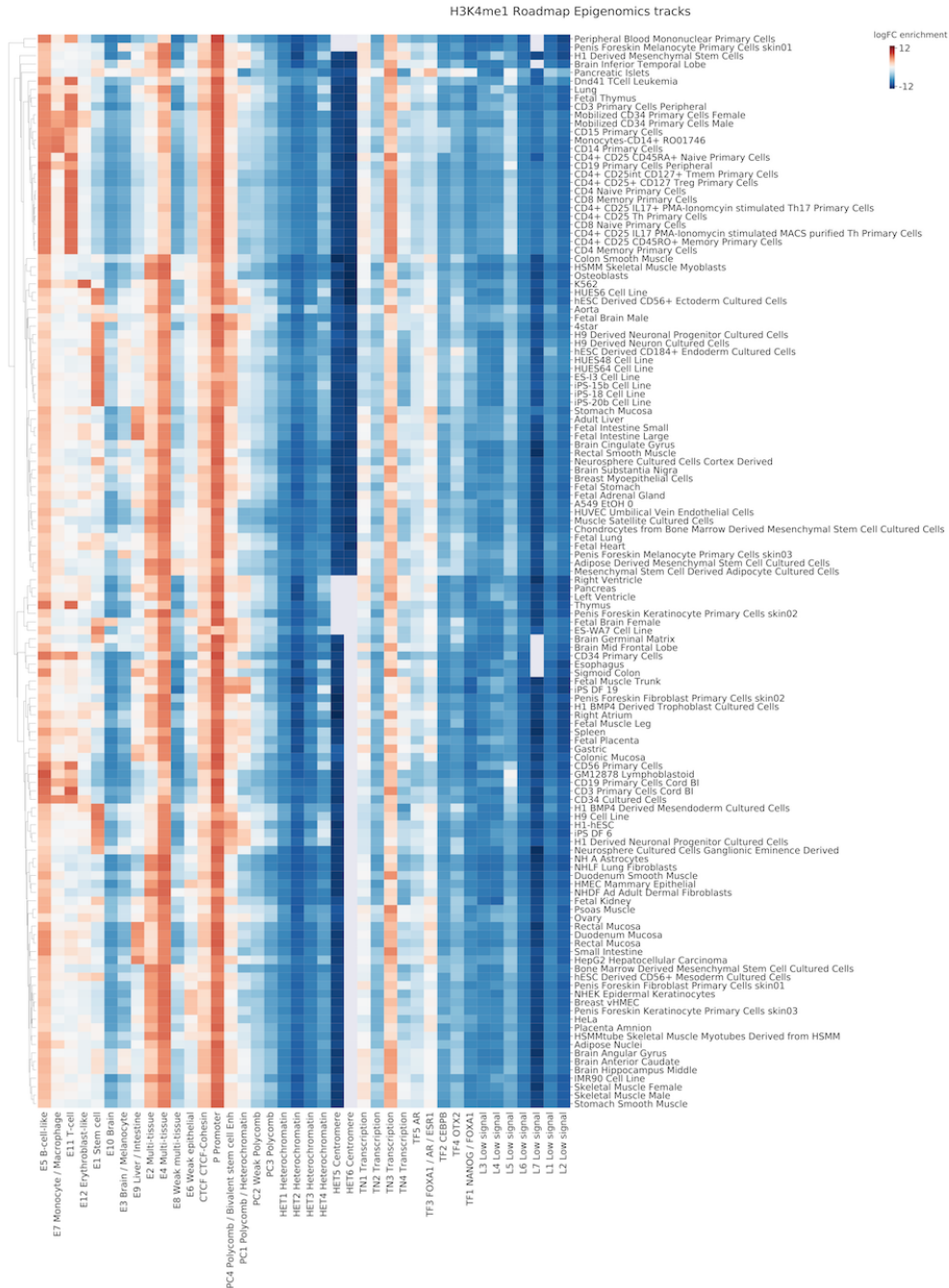
850

851 **Supplementary Figure 7. Enrichment of tissue/cell type-specific H3K4me3 (promoter mark)**

852 **profiles in sequence classes. Log fold change enrichment over genome-average background is shown in**

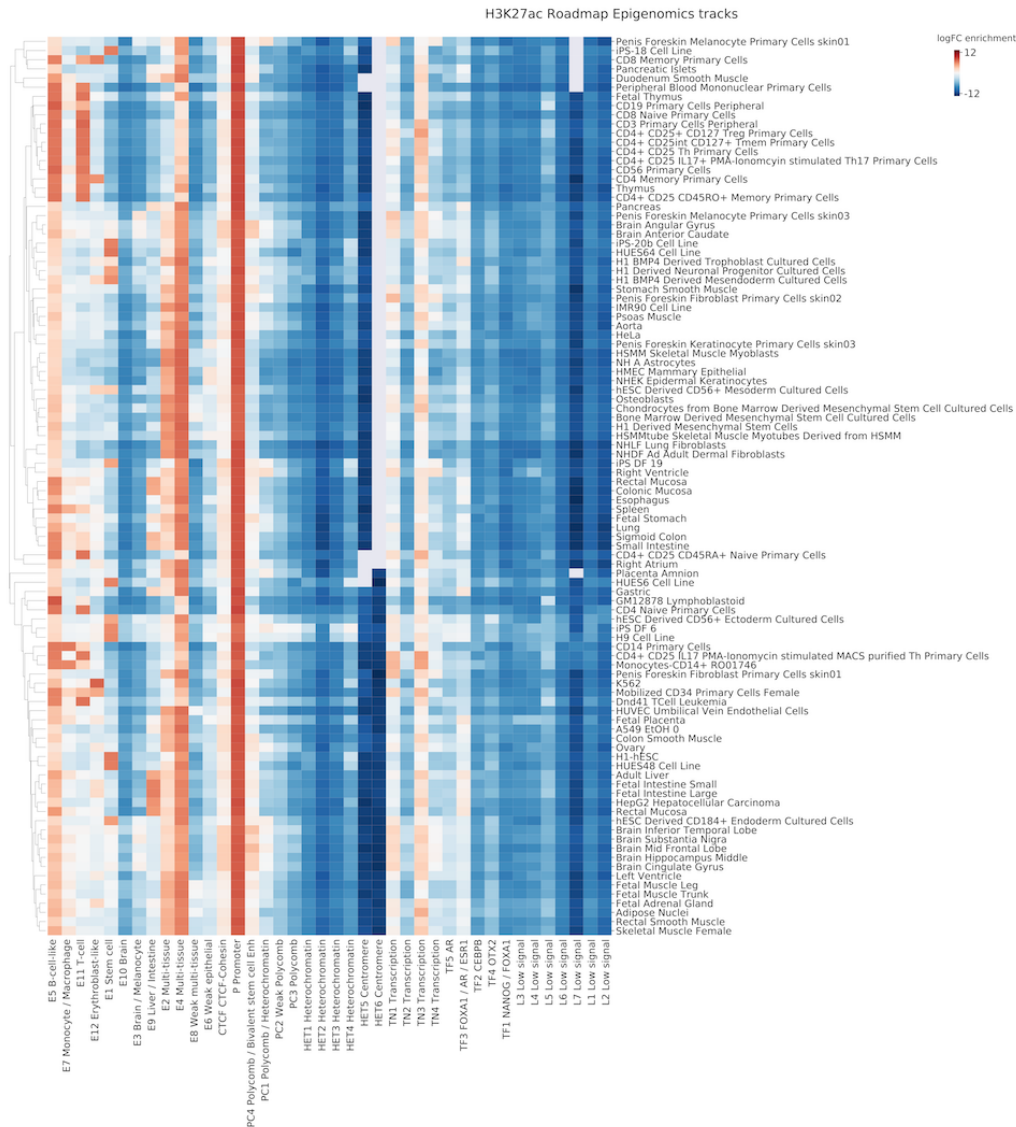
853 **the heatmap. No overlap is indicated by the gray color in the heatmap.**

854

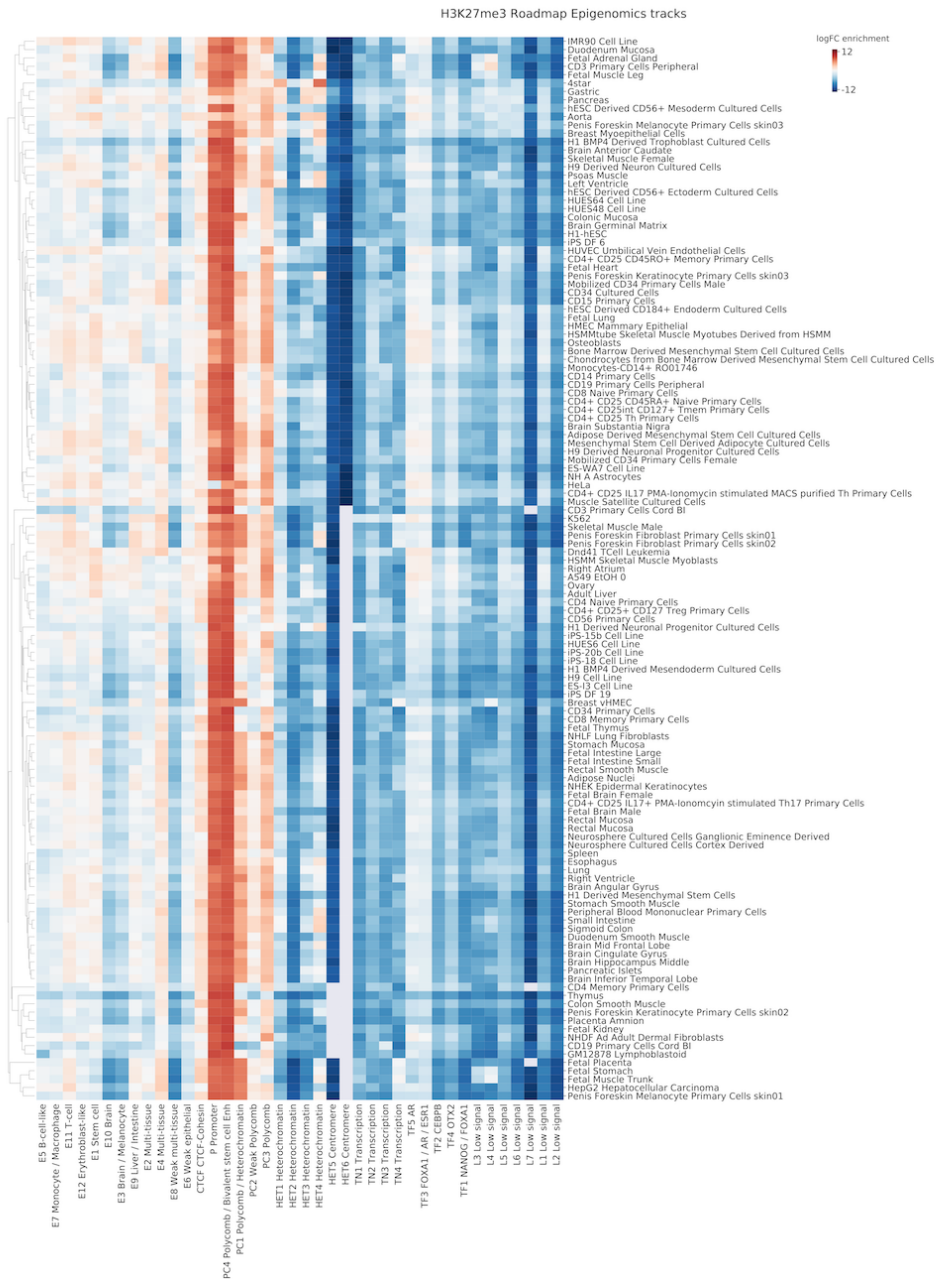


855
856
857
858
859
860

Supplementary Figure 8. Enrichment of tissue/cell type-specific H3K4me1 (enhancer mark) profiles in sequence classes. Log fold change enrichment over genome-average background is shown in the heatmap. No overlap is indicated by the gray color in the heatmap.



861
 862 **Supplementary Figure 9. Enrichment of tissue/cell type-specific H3K27ac (enhancer mark) profiles**
 863 **in sequence classes.** Log fold change enrichment over genome-average background is shown in the
 864 heatmap. No overlap is indicated by the gray color in the heatmap.
 865
 866



867
 868 **Supplementary Figure 10. Enrichment of tissue/cell type-specific H3K27me3 (Polycomb mark)**
 869 **profiles in sequence classes.** Log fold change enrichment over genome-average background is shown in
 870 the heatmap. No overlap is indicated by the gray color in the heatmap.

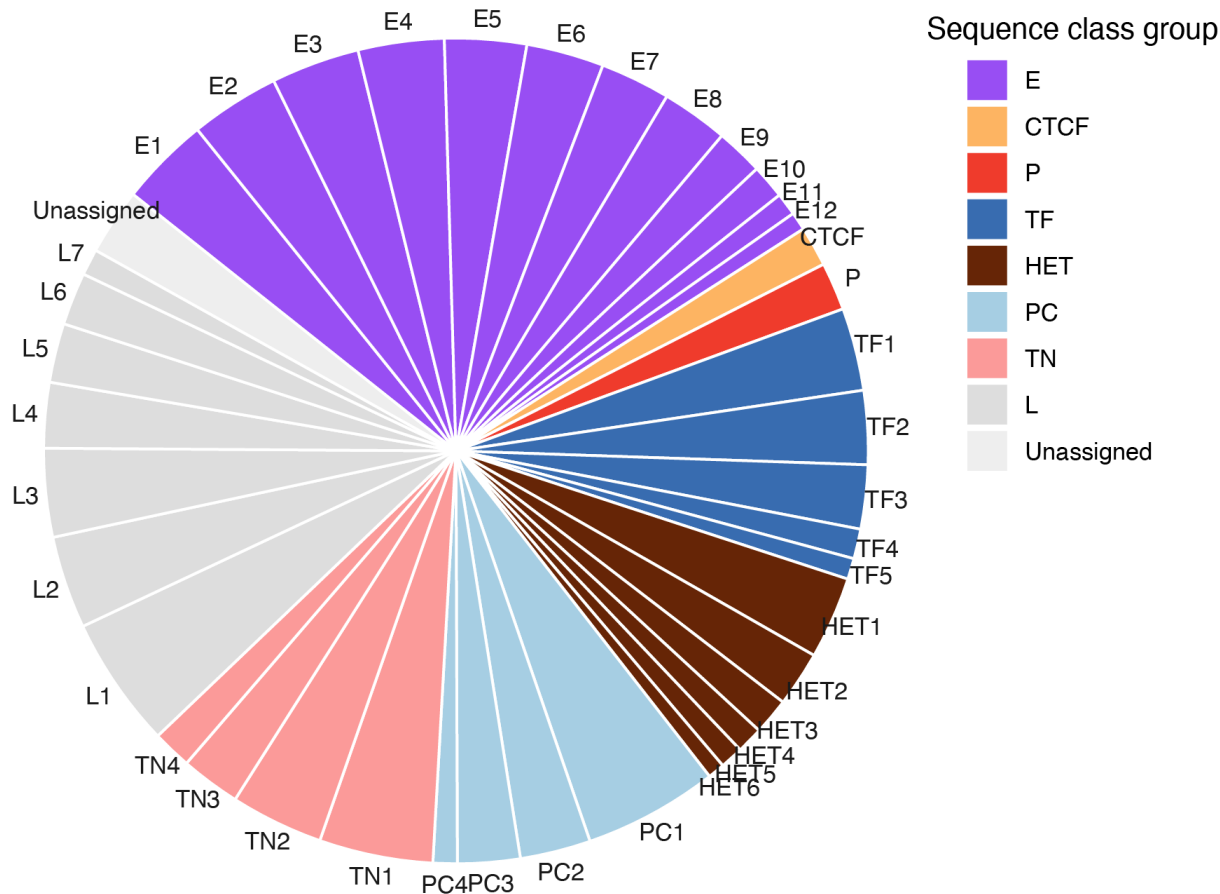
871
 872
 873



879
880
881
882
883
884
885
886

Supplementary Figure 12. Enrichment of tissue/cell type-specific H3K36me3 (transcription mark) profiles in sequence classes. Log fold change enrichment over genome-average background is shown in the heatmap. No overlap is indicated by the gray color in the heatmap.

887



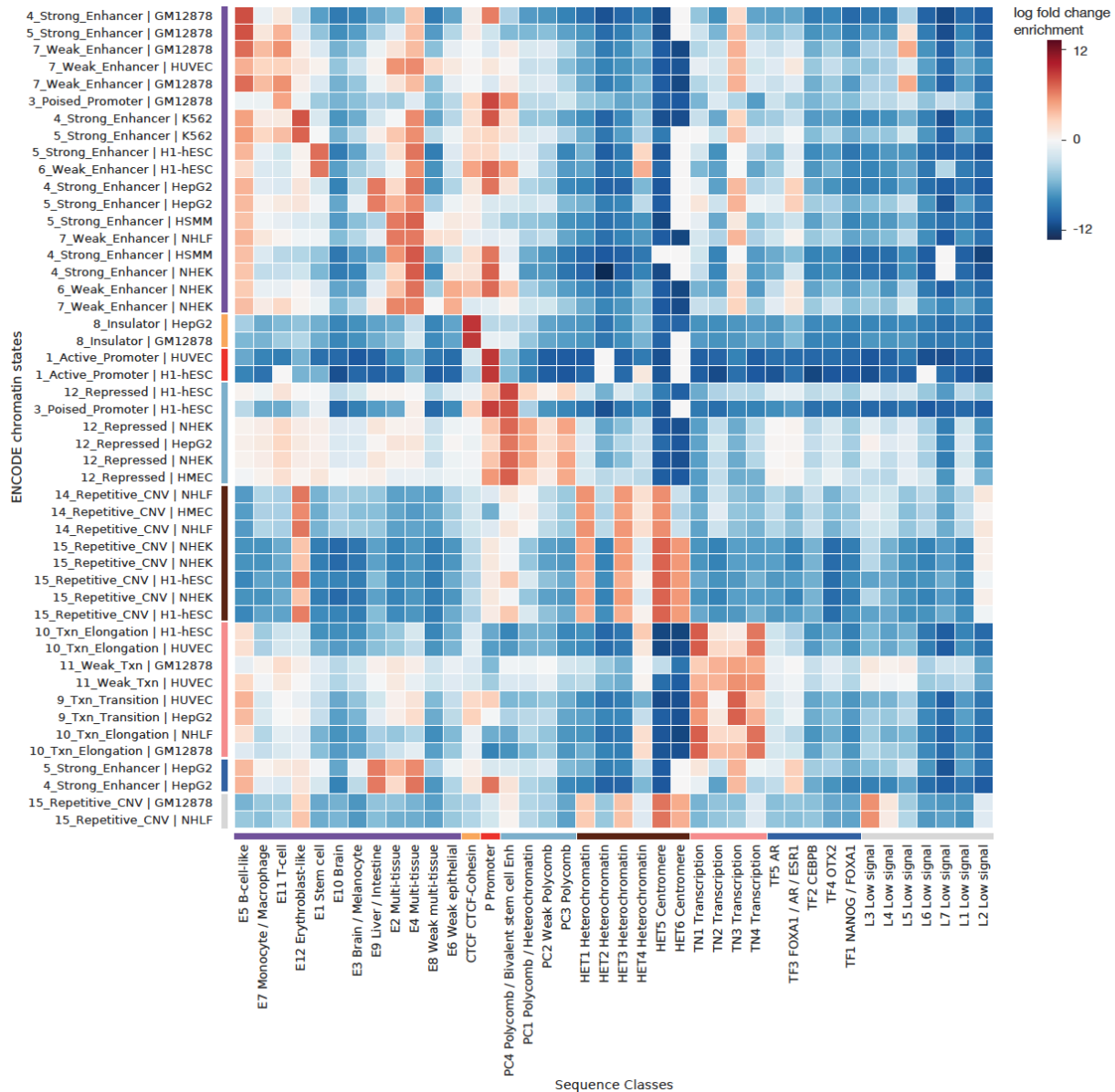
888

889 **Supplementary Figure 13. Genome sequence proportion covered by each sequence class.** The
890 proportion of each sequence class is shown in the pie chart. Genome-wide sequence class assignments
891 were based on Louvain clustering of Sei predictions of sequence tiling the genome with 100bp step size.
892 The clusters unassigned to sequence classes due to the small size (below top 40 clusters) were categorized
893 as “Unassigned”.

894

895

896



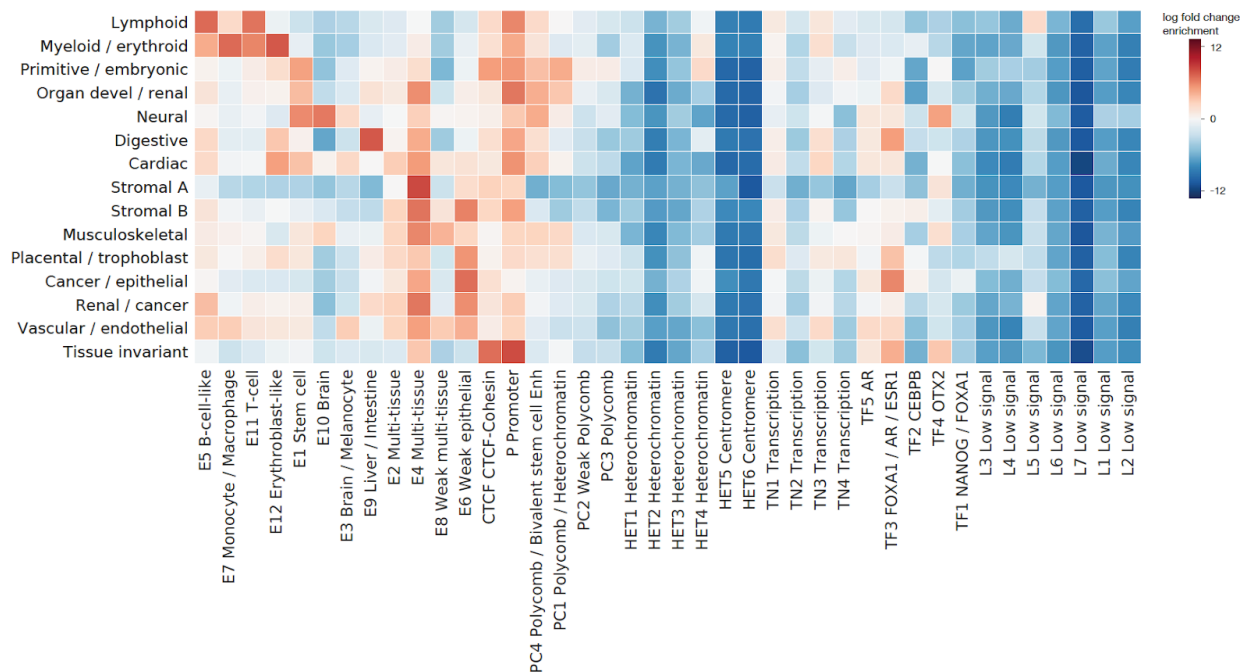
897

898 **Supplementary Figure 14. Sequence-class-specific enrichment of ENCODE chromatin states.** Log

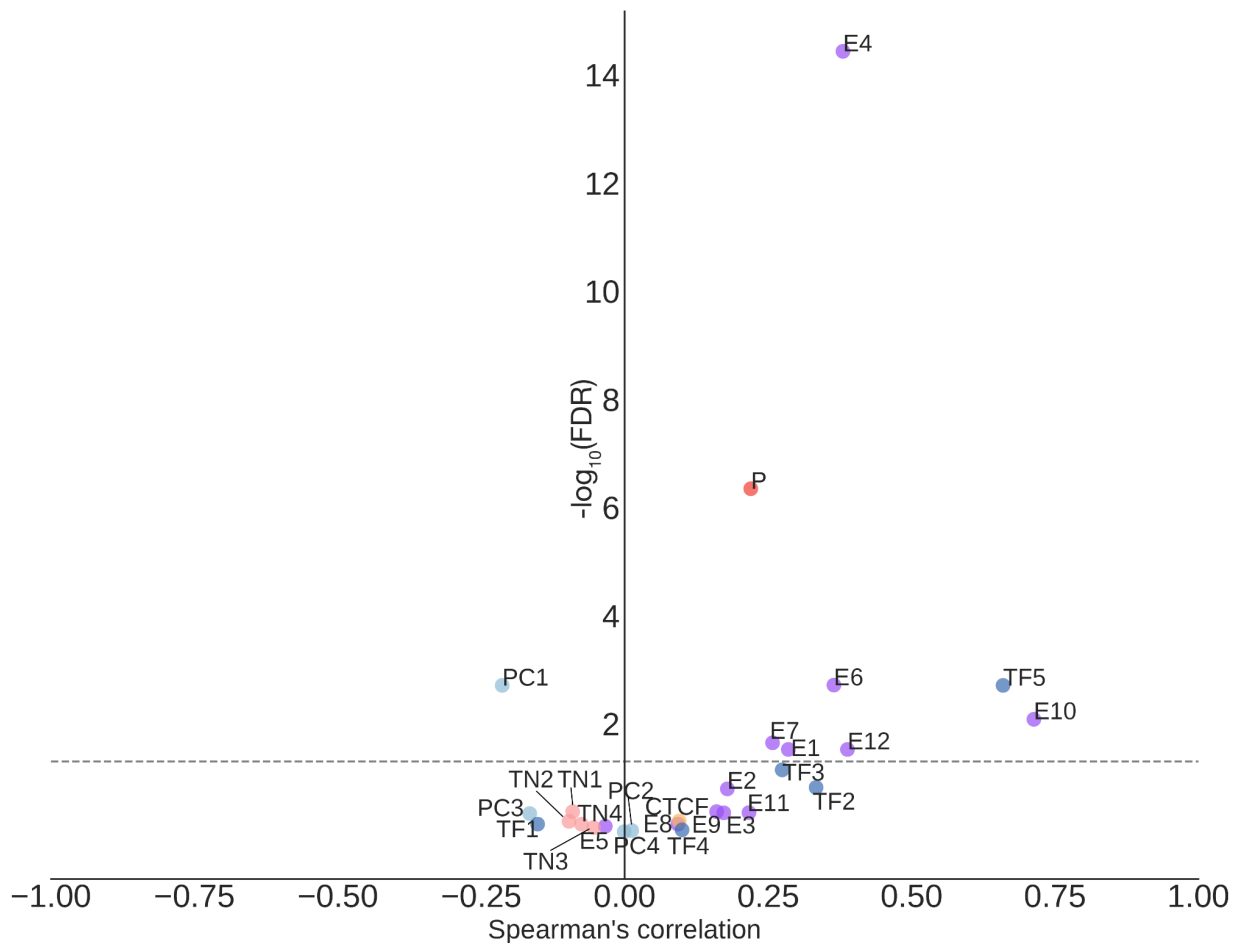
899 fold change enrichment over genome-average background is shown in the heatmap. Top 2 chromatin

900 states enriched were selected for each sequence class.

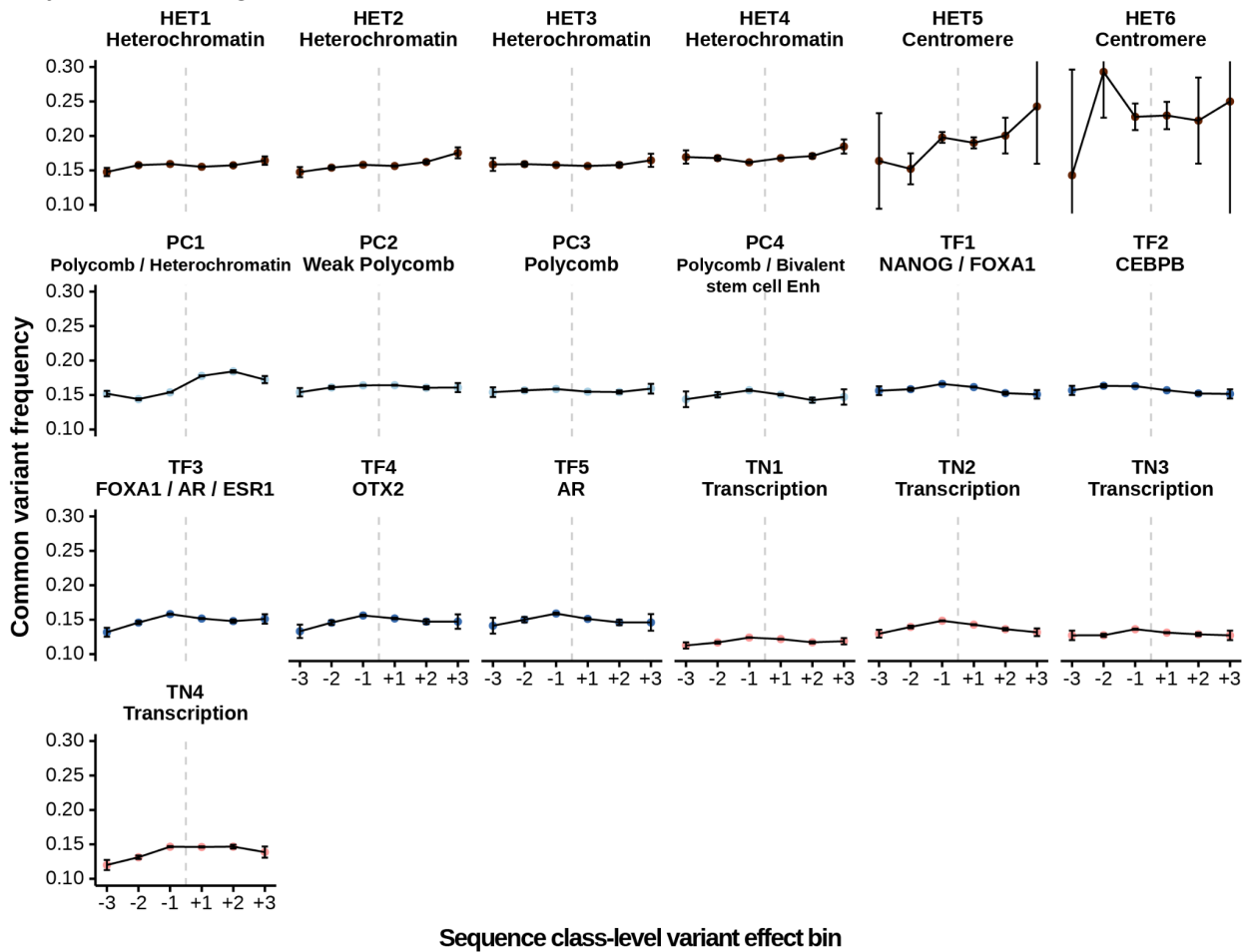
901



902
 903 **Supplementary Figure 15. Sequence-class-specific enrichment of tissue-specific DHS vocabulary**²⁵.
 904 Log fold change enrichment over genome-average background is shown in the heatmap.



906 **Supplementary Figure 16. Regulatory sequence-class-level variant effects for SNPs with PIP > 0.95**
 907 **are predictive of directional GTEx variant gene expression effects.** Variants assigned to sequence
 908 classes based on the sequence class annotation for the reference genome. The x-axis shows Spearman
 909 correlations between the predicted sequence-class-level variant effects and the signed GTEx variant effect
 910 sizes (slopes) and the y-axis shows the corresponding log₁₀ p-values. The dotted gray line denotes the
 911 Benjamini-Hochberg FDR < 0.05 threshold.



912 **Supplementary Figure 17. Population allele frequency profiles for variants in heterochromatin,**
 913 **low signal, polycomb, and transcription sequence classes.** Comparison of common variant frequencies
 914 of 1000 Genomes variants assigned to different sequence classes and variant effect bins. The common
 915 variant threshold is >0.01 allele frequency across the 1000 Genomes population. Error bars show +/- 1
 916 standard error(SE). The sequence-class-level variant effects are assigned to 6 bins (+3: top 1% positive,
 917 +2: top 1%-10% positive, +1, top 10% -100% positive, -3: top 1% negative, -2: top 1%-10% negative, -1,
 918 top 10% -100% negative).
 919

920

947

948 **Supplementary File 5. Significant UKBB GWAS trait - sequence class associations identified with**
949 **LDSR conditioned on the baseline annotations.**

950

951 **Supplementary File 6. Predicted sequence class-level variant effects for HGMD regulatory disease**
952 **mutations.** HGMD regulatory disease mutations with sequence-class level variant effect score >1.1 are
953 included.

954

955 **Supplementary File 7. Detailed Sei model architecture specification.**

956

957 **References**

958 1. Edwards, S. L., Beesley, J., French, J. D. & Dunning, M. Beyond GWASs: Illuminating the
959 dark road from association to function. *American Journal of Human Genetics* vol. 93 779–
960 797 (2013).

961 2. Wieczorek, D. *et al.* A specific mutation in the distant sonic hedgehog (SHH) cis-regulator
962 (ZRS) causes Werner mesomelic syndrome (WMS) while complete ZRS duplications
963 underlie Haas type polysyndactyly and preaxial polydactyly (PPD) with or without
964 triphalangeal thumb. *Hum. Mutat.* **31**, 81–89 (2010).

965 3. Bernstein, B. E. *et al.* An integrated encyclopedia of DNA elements in the human genome.
966 *Nature* **489**, 57–74 (2012).

967 4. Roadmap Epigenomics Consortium *et al.* Integrative analysis of 111 reference human
968 epigenomes. *Nature* **518**, 317–330 (2015).

969 5. Zheng, R. *et al.* Cistrome Data Browser: expanded datasets and new tools for gene
970 regulatory analysis. *Nucleic Acids Res.* **47**, D729–D735 (2019).

971 6. Fillion, G. J. *et al.* Systematic protein location mapping reveals five principal chromatin
972 types in *Drosophila* cells. *Cell* **143**, 212–224 (2010).

973 7. Ernst, J. *et al.* Mapping and analysis of chromatin state dynamics in nine human cell types.
974 *Nature* **473**, 43–49 (2011).

- 975 8. Hoffman, M. M. *et al.* Unsupervised pattern discovery in human chromatin structure
976 through genomic segmentation. *Nat. Methods* **9**, 473–476 (2012).
- 977 9. Zhou, J. & Troyanskaya, O. G. Probabilistic modelling of chromatin code landscape reveals
978 functional diversity of enhancer-like chromatin states. *Nat. Commun.* (2016)
979 doi:10.1038/ncomms10528.
- 980 10. Boix, C. A., James, B. T., Park, Y. P., Meuleman, W. & Kellis, M. Regulatory genomic
981 circuitry of human disease loci by integrative epigenomics. *Nature* **590**, 300–307 (2021).
- 982 11. Alipanahi, B., DeLong, A., Weirauch, M. T. & Frey, B. J. Predicting the sequence
983 specificities of DNA- and RNA-binding proteins by deep learning. *Nat. Biotechnol.* **33**,
984 831–838 (2015).
- 985 12. Zhou, J. & Troyanskaya, O. G. Predicting effects of noncoding variants with deep learning-
986 based sequence model. *Nat. Methods* **12**, 931–934 (2015).
- 987 13. Zhou, J. *et al.* Deep learning sequence-based ab initio prediction of variant effects on
988 expression and disease risk. *Nat. Genet.* (2018) doi:10.1038/s41588-018-0160-6.
- 989 14. Kelley, D. R. *et al.* Sequential regulatory activity prediction across chromosomes with
990 convolutional neural networks. *Genome Res.* (2018) doi:10.1101/gr.227819.117.
- 991 15. Kelley, D. R. Cross-species regulatory sequence activity prediction. *PLoS Comput. Biol.* **16**,
992 e1008050 (2020).
- 993 16. Avsec, Ž. *et al.* Deep learning at base-resolution reveals motif syntax of the cis-regulatory
994 code. *bioRxiv* 737981 (2019) doi:10.1101/737981.
- 995 17. Cofer, E. M. *et al.* Modeling transcriptional regulation of model species with deep learning.
996 *Genome Res.* **31**, 1097–1105 (2021).
- 997 18. Loh, P.-R., Kichaev, G., Gazal, S., Schoech, A. P. & Price, A. L. Mixed-model association

- 998 for biobank-scale datasets. *Nat. Genet.* **50**, 906–908 (2018).
- 999 19. McInnes, L., Healy, J., Saul, N. & Großberger, L. UMAP: Uniform Manifold
1000 Approximation and Projection. *Journal of Open Source Software* (2018)
1001 doi:10.21105/joss.00861.
- 1002 20. Poličar, P. G., Stražar, M. & Zupan, B. openTSNE: a modular Python library for t-SNE
1003 dimensionality reduction and embedding. *BioRxiv* (2019).
- 1004 21. Blondel, V. D., Guillaume, J.-L., Lambiotte, R. & Lefebvre, E. Fast unfolding of
1005 communities in large networks. *J. Stat. Mech.* **2008**, P10008 (2008).
- 1006 22. Hawkins, R. D. *et al.* Distinct epigenomic landscapes of pluripotent and lineage-committed
1007 human cells. *Cell Stem Cell* **6**, 479–491 (2010).
- 1008 23. Boros, J., Arnoult, N., Stroobant, V., Collet, J.-F. & Decottignies, A. Polycomb repressive
1009 complex 2 and H3K27me3 cooperate with H3K9 methylation to maintain heterochromatin
1010 protein 1 α at chromatin. *Mol. Cell. Biol.* **34**, 3662–3674 (2014).
- 1011 24. Schwämmle, V. *et al.* Systems Level Analysis of Histone H3 Post-translational
1012 Modifications (PTMs) Reveals Features of PTM Crosstalk in Chromatin Regulation. *Mol.*
1013 *Cell. Proteomics* **15**, 2715–2729 (2016).
- 1014 25. Meuleman, W. *et al.* Index and biological spectrum of human DNase I hypersensitive sites.
1015 *Nature* **584**, 244–251 (2020).
- 1016 26. The GTEx Consortium. The GTEx Consortium atlas of genetic regulatory effects across
1017 human tissues. *Science* **369**, 1318–1330 (2020).
- 1018 27. Kerimov, N. *et al.* A compendium of uniformly processed human gene expression and
1019 splicing quantitative trait loci. *Nat. Genet.* **53**, 1290–1299 (2021).
- 1020 28. 1000 Genomes Project Consortium, T. 1000 G. P. *et al.* An integrated map of genetic

- 1021 variation from 1,092 human genomes. *Nature* (2012) doi:10.1038/nature11632.
- 1022 29. Finucane, H. K. *et al.* Partitioning heritability by functional annotation using genome-wide
1023 association summary statistics. *Nat. Genet.* **47**, 1228–1235 (2015).
- 1024 30. Finucane, H. K. *et al.* Heritability enrichment of specifically expressed genes identifies
1025 disease-relevant tissues and cell types. *Nat. Genet.* **50**, 621–629 (2018).
- 1026 31. Reshef, Y. A. *et al.* Detecting genome-wide directional effects of transcription factor
1027 binding on polygenic disease risk. doi:10.1101/204685.
- 1028 32. Sudlow, C. *et al.* UK biobank: an open access resource for identifying the causes of a wide
1029 range of complex diseases of middle and old age. *PLoS Med.* **12**, e1001779 (2015).
- 1030 33. Paththinige, C. S., Sirisena, N. D. & Dissanayake, V. H. W. Genetic determinants of
1031 inherited susceptibility to hypercholesterolemia – a comprehensive literature review. *Lipids*
1032 *Health Dis.* **16**, 103 (2017).
- 1033 34. Stenson, P. D. *et al.* The Human Gene Mutation Database: 2008 update. *Genome Med.* **1**, 13
1034 (2009).
- 1035 35. Al-Qattan, M. M., Al Abdulkareem, I., Al Haidan, Y. & Al Balwi, M. A novel mutation in
1036 the SHH long-range regulator (ZRS) is associated with preaxial polydactyly, triphalangeal
1037 thumb, and severe radial ray deficiency. *Am. J. Med. Genet. A* **158A**, 2610–2615 (2012).
- 1038 36. Gurnett, C. A. *et al.* Two novel point mutations in the long-range SHH enhancer in three
1039 families with triphalangeal thumb and preaxial polydactyly. *Am. J. Med. Genet. A* **143A**,
1040 27–32 (2007).
- 1041 37. De Castro-Orós, I. *et al.* Functional analysis of LDLR promoter and 5' UTR mutations in
1042 subjects with clinical diagnosis of familial hypercholesterolemia. *Hum. Mutat.* **32**, 868–872
1043 (2011).

- 1044 38. Ilkovski, B. *et al.* Mutations in PIGY: expanding the phenotype of inherited
1045 glycosylphosphatidylinositol deficiencies. *Hum. Mol. Genet.* **24**, 6146–6159 (2015).
- 1046 39. Sakai, T., Ohtani, N., McGee, T. L., Robbins, P. D. & Dryja, T. P. Oncogenic germ-line
1047 mutations in Sp1 and ATF sites in the human retinoblastoma gene. *Nature* **353**, 83–86
1048 (1991).
- 1049 40. Plenge, R. M. *et al.* A promoter mutation in the XIST gene in two unrelated families with
1050 skewed X-chromosome inactivation. *Nat. Genet.* **17**, 353–356 (1997).
- 1051 41. Pugacheva, E. M. *et al.* Familial cases of point mutations in the XIST promoter reveal a
1052 correlation between CTCF binding and pre-emptive choices of X chromosome inactivation.
1053 *Hum. Mol. Genet.* **14**, 953–965 (2005).
- 1054 42. De Gobbi, M. *et al.* A regulatory SNP causes a human genetic disease by creating a new
1055 transcriptional promoter. *Science* (2006) doi:10.1126/science.1126431.
- 1056 43. Horn, S. *et al.* TERT promoter mutations in familial and sporadic melanoma. *Science* **339**,
1057 959–961 (2013).
- 1058 44. Surrey, S., Delgrosso, K., Malladi, P. & Schwartz, E. A single-base change at position -175
1059 in the 5'-flanking region of the G gamma-globin gene from a black with G gamma-beta+
1060 HPFH. *Blood* **71**, 807–810 (1988).
- 1061 45. Amemiya, H. M., Kundaje, A. & Boyle, A. P. The ENCODE Blacklist: Identification of
1062 Problematic Regions of the Genome. *Sci. Rep.* **9**, 9354 (2019).
- 1063 46. Chen, K. M., Cofer, E. M., Zhou, J. & Troyanskaya, O. G. Selene: a PyTorch-based deep
1064 learning library for sequence data. *Nat. Methods* (2019) doi:10.1038/s41592-019-0360-8.
- 1065