1

**Universal annotation of the human genome through integration of over a thousand epigenomic**

**datasets**

**Ha Vu[1,2], Jason Ernst[1,2,3,4,5,6,7]**

[1] Bioinformatics Interdepartmental Program, University of California, Los Angeles, CA, 90095, USA.

[2] Department of Biological Chemistry, University of California, Los Angeles, Los Angeles, CA 90095, USA

[3] Eli and Edythe Broad Center of Regenerative Medicine and Stem Cell Research at University of California, Los Angeles, Los Angeles, CA 90095, USA

[4] Computer Science Department, University of California, Los Angeles, Los Angeles, CA 90095, USA

[5] Jonsson Comprehensive Cancer Center, University of California, Los Angeles, Los Angeles, CA 90095, USA

[6] Molecular Biology Institute, University of California, Los Angeles, Los Angeles, CA 90095, USA

[7] Department of Computational Medicine, University of California, Los Angeles, Los Angeles, CA 90095, USA

16

17

18   **Abstract**

19   Genome-wide maps of chromatin marks such as histone modifications and open chromatin sites

20   provide valuable information for annotating the non-coding genome, including identifying

21   regulatory elements. Computational approaches such as ChromHMM have been applied to

22   discover and annotate chromatin states defined by combinatorial and spatial patterns of

23   chromatin marks within the same cell type. An alternative 'stacked modeling' approach was

24   previously suggested, where chromatin states are defined jointly from datasets of multiple cell

25   types to produce a single universal genome annotation based on all datasets. Despite its potential

26   benefits for applications that are not specific to one cell type, such an approach was previously

27   applied only for small-scale specialized purposes. Large-scale applications of stacked modeling

28   have previously posed scalability challenges. In this paper, using a version of ChromHMM

29   enhanced for large-scale applications, we applied the stacked modeling approach to produce a

30   universal chromatin state annotation of the human genome using over 1000 datasets from more

31   than 100 cell types, denoted the full-stack model. The full-stack model states show distinct

32   enrichments for external genomic annotations, which we used in characterizing each state.

33   Compared to cell-type-specific annotations, the full-stack annotation directly differentiates

34   constitutive from cell-type-specific activity and is more predictive of locations of external

35   genomic annotations. Overall, the full-stack ChromHMM model provides a universal chromatin

36   state annotation of the genome and a unified global view of over 1000 datasets. We expect this to

37   be a useful resource that complements existing cell-type-specific annotations for studying the

38   non-coding human genome.

39

**Introduction**

Genome-wide maps of histone modifications, histone variants and open chromatin provide valuable information for annotating the non-coding genome features, including various types of regulatory elements [1–5]. These maps -- produced by assays such as chromatin immunoprecipitation followed by high-throughput sequencing to map histone modifications or DNase-seq to map open chromatin-- can facilitate our understanding of regulatory elements and genetic variants that are associated with disease [6–12]. Efforts by large scale consortia as well as many individual labs have resulted in these maps for many different human cell and tissue types for multiple different chromatin marks [1,9,13–20].

The availability of maps for multiple different chromatin marks in the same cell or tissue type motivated the development of methods such as ChromHMM and Segway that learn 'chromatin states' based on the combinatorial and spatial patterns of marks in such data [21–23]. These methods then annotate genomes in a cell-type-specific manner based on the learned chromatin states. They have been applied to annotate more than a hundred diverse cell and tissue types [3,16,24]. Previously, large collections of cell-type-specific annotations have been generated using either (1) independent models that learn a different set of states in each cell or tissue type or (2) a model that is learned across all cells and tissues, resulting in a common set of states across cell types, yet generating cell-type-specific state annotations. This latter approach has previously been referred to as a 'concatenated' approach (**Supp. Fig. 1**) [22,25]. Variants of this approach allow information from other cell types to influence the state annotations in one cell type at a position, but still produce cell-type-specific state annotations [26,27]. These models that produce cell-type-specific annotations are natural for cell-type-specific analyses.

A complementary approach to applying ChromHMM to data across multiple different cell types referred to as the 'stacked' modeling approach was also previously suggested (**Supp. Fig. 1**) [22,25]. Instead of learning cell-type-specific annotations based on a limited number of datasets per cell type, the stacked modeling approach learns a single universal genome annotation based on the combinatorial and

3

65   spatial patterns in datasets from multiple marks across multiple cell types. This approach differs from the

66   concatenated and independent modeling approaches as those approaches only identify combinatorial and

67   spatial patterns present among datasets within one cell type.

68         Such a universal annotation from stacked modeling provides potential complementary benefits to

69   existing cell-type-specific chromatin state annotations. First, stacked models may help differentiate regions

70   with constitutive chromatin activities from those with cell-type-specific activities. Previously, specific

71   chromatin states from 'concatenated' cell-type-specific annotations were post-hoc clustered to analyze

72   chromatin dynamics across cell and tissue types, yet such an approach does not provide a systematic and

73   global view of the dynamics of all the data [3,16]. Second, the stacked modeling approach bypasses the

74   need to pick a specific cell or tissue type when analyzing a single partitioning and annotation of the genome.

75   Focusing on a single cell or tissue type may not be desirable for many analyses involving other annotations

76   that are not inherently cell-type-specific, such as those involving conserved DNA sequence or genetic

77   variants. Alternatively, compared to analyzing chromatin state annotations across all cells or tissue types,

78   while the stacked model state definitions are more complex, the resulting genome annotations are simpler

79   and non-overlapping. With the stacked modeling, each location is simply assigned to one of $N$ universal

80   states, whereas in the concatenated model, each location is assigned to one of $M$ states in $K$ cell types. The

81   value of $N$ can be selected to be much smaller than the number of possible combinations of chromatin state

82   annotations across cell types at a location with the concatenated modeling, $M^K$, as well as the number of

83   possible combinations of cell types and states, $M*K$. Finally, annotations by the stacked modeling leverages

84   a larger set of data for annotation, and thus has the potential to be able to identify genomic elements with

85   greater sensitivity and specificity.

86         Despite the potential complementary advantages of the 'stacked' modeling approach, it has only

87   been applied on a limited scale to combine data from a small number of cell types for highly specific

88   purposes [28,29]. No large-scale application of the stacked modeling approach to many diverse cell and

89   tissue types has been previously demonstrated. This may have in part been due to large-scale applications

90   of stacked modeling raising scalability challenges not present in cell-type-specific modeling.

4

91    Here, we present a large-scale application of the stacked modeling approach with more than a

92 thousand human epigenomic datasets as input, using a version of ChromHMM of which we enhanced the

93 scalability. We conduct various enrichment analyses on the states resulting from the stacked modeling and

94 give biological interpretations to them. We show that compared to the cell-type-specific annotations, the

95 stacked model's annotation shows greater correspondence to various external genomic annotations not used

96 in the model learning. We analyze the states in terms of enrichment with different types of genetic variants,

97 and highlight specific states of the stacked model that are enriched with phenotypically associated genetic

98 variants. Additionally, we identify specific states enriched with cancer-associated somatic mutations. We

99 expect the stacked model annotations and detailed characterization of the states that we provide will be a

100 valuable resource for studying the epigenome and non-coding genome, complementing existing cell-type-

101 specific annotations.

102

103     **Results**

104     *Annotating the human genome into universal chromatin states*

105     We used the stacked modeling approach of ChromHMM to produce a universal chromatin state

106     annotation of the genome based on data from over 100 cell and tissue types from the Roadmap Epigenomics

107     and ENCODE projects (**Fig. 1**) [14,16]. In total we applied ChromHMM to 1032 datasets for 30 histone

108     modifications, a histone variant (H2A.Z), and DNase I hypersensitivity (**Supp. Fig. 2**). The set of cell and

109     tissue types were the same as those for which cell-type-specific annotations were previously generated by

110     applying the 'concatenated' modeling approach of ChromHMM [22,25]. We note that not all chromatin

111     marks were profiled in all cell or tissue types, but the stacked modelling can still be applied directly.

112     We focused our analysis on a model with 100 states. We used a larger number of states than

113     typically used for cell-type-specific models to reflect the greater information available when defining states

114     based on data from many cell types. At the same time, we limited the model to 100 states to ensure

115     manageable biological interpretation of different states (**Supp. Fig. 3**) (**Methods**). We denote the model's

116     output chromatin state annotation the 'full-stack' genome annotation.

117

118     *Major groups of full-stack states*

119     We characterized each state of the model by analyzing the model parameters (emission probabilities

120     and transition probabilities) and state enrichments for other genome annotations (**Fig.2, 3A, Supp. Fig. 4-**

121     **7**). The other genomic annotations include previous cell-type-specific chromatin state annotations (**Supp.**

122     **Fig. 8**), cell-type-specific gene expression data (**Supp. Fig. 9-10**), and various independent existing

123     genomic annotations (**Fig. 3A**). These independent genomic annotations included annotated gene features,

124     evolutionary constrained elements, and assembly gaps, among others (**Methods**).

125     These analyses led us to group the 100 full-stack states into 16 groups (**Fig. 2A**). One group

126     includes states associated with assembly gaps (state GapArtf1) and alignment artifacts (states GapArtf2-3).

127     Some other groups are associated with repressive or inactive states, including quiescent states (states

128     Quies1-5) (low emissions of all experiments, except possibly weak signals in H3K9me3), heterochromatin

6

129    states associated with H3K9me3 (states HET1-9), and polycomb repressed states associated with

130    H3K27me3 (states ReprPC1-9). There is an acetylations group marked primarily by high emission of

131    various acetylation marks (states Acet1-8). We also formed active and weak candidate enhancers groups

132    (states EnhW1-8 and EnhA1-20, respectively) associated with H3K4me1, DNase,  H2A.Z, and/or

133    H3K27ac. Four groups are associated with transcriptional activities, including a group of transcribed

134    enhancers (states TxEnh1-8), two groups of weak or strong transcription (states TxWk1-2, Tx1-8,

135    respectively), and one group associated with exon and transcription (states TxEx1-4). These transcriptional

136    activities groups are associated with at least one of these marks H3K36me3, H3K79me1, H3K79me2, and

137    H4K20me1. Another group consists of two zinc finger (ZNF) gene states associated with H3K36me3 and

138    H3K9me3 (states ZNF1-2). A DNase group consists of one state (DNase1) with strong emission of only

139    DNase in all profiled cell types.  Three groups are associated with promoter activities, marked by emission

140    of some promoter marks such as H3K4me3, H3K4me2, and H3K9ac. One promoter group was of bivalent

141    states associated with promoter marks and H3K27me3 (states BivProm1-4). The other two promoter groups

142    were flanking promoter states (PromF1-7) and transcription start sites (TSS) states (TSS1-2) where the

143    flanking promoter states also show emission of H3K4me1.

144        Enrichments for external annotations supported these state groupings (**Fig. 3A**), as well as further

145    distinctions or characterizations among states within each group. For example, the state in the assembly gap

146    group (GapArt1) had ~8 fold enrichment for assembly gaps and contained 99.99% of all assembly gaps

147    (**Fig. 3A**). The states in the zinc finger gene group, ZNF1-2, had 20.8 and 68.6 fold enrichment for zinc

148    finger named genes, respectively (**Fig. 3A**). States in the transcription groups (TxEnh1-8, TxWk1-2, Tx1-

149    8, TxEx1-4) were all at least 2.1 fold enriched for annotated genes, which covered 88.8–97.5% of the states.

150    These states are associated with higher expression of genes across different cell types, particularly when

151    downstream of their TSS (**Fig 3A, C, Supp. Fig. 9-10**). Distinctions were seen among these states, for

152    example, in terms of their positional enrichments relative to TES (**Fig. 3A, D, Supp. Fig. 11**). States in the

153    flanking promoter group (PromF1-7) showed 6.5-28 fold enrichment for being within 2kb of annotated

154    TSS, and genes whose TSS regions overlapped these states had higher average gene expression across

7

155 different cell types (**Fig. 3A, C, Supp. Fig. 9-10**). These states differed among each other in their

156 enrichments with upstream or downstream regions of the TSS (**Fig. 3E, Supp. Fig. 11**). The states in the

157 transcription start site group (TSS1-2) had particularly high enrichments around the TSS, with >= 100 fold

158 enrichment (**Fig. 3A, E**). The DNase specific group, which comprised of the state DNase1, showed strong

159 enrichment for CTCF-specific chromatin states defined in six cell types [30] (**Fig. 3F, Supp. Fig. 12**). A

160 detailed characterization of all states can be found in **Supplementary Data**.

161

162 *Stacked Model Differentiates Cell-Type-Specific from Constitutive Activity*

163 While the major groups of states outlined above can correspond to states in cell-type-specific models [3,16],

164 the full-stack states provide additional information. For example, the states directly differentiate cell-type-

165 specific from constitutive activities. Consistent with previous findings that enhancers tend to be relatively

166 cell-type-specific while promoters tend to be shared across cell types [3,31], enhancer states exhibited

167 clearer cell-type-specific associations than those of the promoter states (**Figure 2C, Supplementary**

168 **Data.**). This is also reflected in the states' coefficients of variation across different cell groups in terms of

169 emission probabilities for the marks DNase, H3K27ac, H3K4me1, H3K4me2, H3K4me3 and H3K9ac. On

170 average, states of enhancer and weak enhancer groups (EnhW1-8, EnhA1-20) show at least two fold higher

171 of the coefficients of variations compared to states in the TSS, flanking promoter and bivalent promoters

172 groups (TSS1-2, PromF1-7, BivProm4) (**Supp. Fig. 13**). The enhancer states differed among each other in

173 their associations with different cell/tissue types such as brain (EnhA6), blood (EnhA7-9 and EnhWk6),

174 digestive tissue (EnhA14-15), and embryonic stem cells (EnhA18) (**Fig. 1-2, Supp. Fig. 14-15**). These

175 differences in cell-type-specific activities are also associated with different gene expression levels of

176 overlapping genes with the states. For example, some blood enhancer states (EnhA8, EnhA9, EnhWk6)

177 overlapped genes with higher average gene expression in cell types of the blood group, while some enhancer

178 states specific to digestive group or liver tissues (EnhA14, EnhA15) showed higher gene expression in the

179 corresponding cell or tissue types (**Fig. 3C, Supp. Fig. 9**).

180    Other groups of states besides enhancers also had individual states with cell-type-specific

181    differences. For example, four of the nine states in the heterochromatin group show higher emission

182    probabilities of H3K9me3 in only subsets of cell types (states HET1-2 with IMR90 and Epithelial cells;

183    state HET4 with adipose, mesench, neurospheres, ESC, HSC&B-cells; state HET9 with ESC/iPSC)

184    (**Supplementary Data**). In addition, some quiescent states (Quies1-2, Quies4-5) show weak signals of

185    H3K9me3 in specific groups of cell types (**Supplementary Data**). States in the polycomb repressed and

186    bivalent promoter groups (ReprPC1-9, BivProm1-4) also show differences in signals across cell groups,

187    such as state ReprPC9, which showed H3K27me3 signals in only ESC/iPSC cell types (**Supplementary**

188    **Data**). The ability of the stacked modeling approach to explicitly annotate both cell-type-specific and

189    constitutive patterns for diverse classes of chromatin states highlights an advantage of this approach relative

190    to the concatenated modeling.

191

192    *Full-stack states are more predictive of external annotations than cell-type-specific models*

193    Another benefit of the stacked modeling approach is its ability to more accurately identify genomic

194    elements shared across cell and tissue types. To demonstrate this, we compared the full-stack state

195    annotation against two sets of cell-type-specific chromatin state annotations in terms of recovering locations

196    of various external genome annotations (**Methods**). One set was the previously published 18-chromatin

197    state annotations defined in 98 cell or tissue types (equivalently, reference epigenomes) using a common

198    set of six chromatin marks from Roadmap Epigenomics with the concatenated modeling approach [16].

199    The other set of annotations we compared the full-stack annotation against were 100-state cell-type-specific

200    annotations that we generated separately for each of the 127 cell or tissue types using all available chromatin

201    marks in the respective cell or tissue type (**Methods**).

202    The external genome annotations we used for the evaluations included locations of coding

203    sequences, assembly gaps, CpG Islands, lamina associated domains (laminB1lads), PhastCons elements,

204    pseudogenes, exons, gene body, transcription end sites, transcription start sites and the 2kb neighboring

205    regions, repeat elements, and zinc finger named (ZNF) genes. The full-stack annotation resulted in the best

9

206    AUROC (area under the receiver operating curve) in predicting all genomic annotations compared to the

207    previous 18-state cell-type-specific annotations across all cell types (**Supp. Fig. 16-17**). The full-stack

208    model also showed the best AUROC in recovering locations of these genomic annotations for 10 of 13

209    evaluations compared to all 100-state cell-type-specific annotations (**Fig. 4**). The only evaluations in which

210    the full-stack model did not outperform all 100-state cell-type-specific models were those involving

211    assembly gaps, laminB1lads, and ZNF genes (**Fig. 4B, Supp. Fig. 18**), where at most 6 of the 127 100-state

212    cell-type-specific models performed better. Additionally, we obtained similar results in comparing full-

213    stack annotations with cell-type-specific annotations in predicting CTCF specific chromatin states in

214    multiple cell types, where the full-stack annotation resulted in highest AUROC in all cases (**Supp. Fig.**

215    **19**).

216          Overall, these results demonstrate the benefits of full-stack chromatin state annotations, which

217    showed better predictive power in recovering the locations of a variety of independent genomic annotations.

218    The increased predictive power of the stacked modeling approach can be attributed to it taking into account

219    information from more datasets that cover a large number of cell types when inferring state annotations.

220

221    *Full-stack states show distinct enrichments for repeat elements*

222    As the full-stack model showed greater predictive power for repeat elements than cell-type-specific models

223    (**Fig. 4A, Supp. Fig. 29-31**), we next analyzed which states contributed most to this power. The full-stack

224    state enrichments for bases in repeat elements ranged from 10-fold depletion to 2-fold enrichment (**Fig.**

225    **3A**). The top ten states most enriched with repeat elements were chromatin states that were associated with

226    H3K9me3 marks and in the heterochromatin, artifact, quiescent, or ZNF genes groups (**Fig. 5A-B**).

227    We also observed that individual full-stack states had distinct enrichments for different repeat classes (**Fig.**

228    **5C, Supp. Fig. 20**). For example, Acet1, a state associated with various acetylation marks and H3K9me3

229    had a 23-fold enrichment for simple repeats (**Supp. Fig. 20**). The two states in the artifact group, GapArtf2-

230    3, had a particularly high enrichment for satellite (181 and 145 fold, respectively) and rRNA repeat classes

231    (75 and 580 fold, respectively) (**Fig. 5C, Supp. Fig. 20**). States in the transcription start site group, TSS1-

10

232    2, were most strongly enriched with tRNA and low complexity repeat class (~10-60 fold) (**Supp. Fig. 20**).

233    We also saw specific states associated with the largest repeat classes of the genome, SINEs, LINEs, and

234    LTRs. SINE repeats were most enriched in state Tx5 (3.7 fold) (**Fig. 5C**), which had high emission of

235    H3K36me3 (**Fig. 2A-B, Supp. Fig. 4-5**). LINEs and LTRs repeats were most enriched for states in the

236    H3K9me3-associated heterochromatin group with LINE most enriched in HET3 (3.4 fold), while LTRs

237    were most enriched in HET5 (4.7 fold) (**Fig. 5C, Supp. Fig. 20**). We also confirmed that the increased

238    predictive power of the full-stack model over cell-type-specific models, which was previously seen for

239    repeat elements overall, also held for most of the individual repeat classes (**Supp. Fig. 21**).

240

241    *Full-stack states show distinct enrichments for constrained elements and conservation states*

242    Sequence constrained elements are another class of genomic elements that are not cell-type-specific and for

243    which the full-stack model showed greater predictive power than the cell-type-specific models (**Fig. 4B,**

244    **Supp. Fig. 16-18**). We next sought to better understand the relationship between full-stack states and

245    sequence conservation annotations. We observed 10 states that had at least a 3.4 fold enrichment for

246    PhastCons elements (**Fig. 5A**). These states were associated with the TSSs or being proximal to them

247    (TSS1-2 and PromF4-5), transcription with strong H3K36me3 signals (TxEx2 and TxEnh4), or enhancers

248    associated with mesenchymal, muscle, heart, neurosph, adipose (EnhA2) (**Fig. 5A-B**). In contrast, seven

249    states (HET3-4,6-7,9, Quies4, Gap Artf2) were more than two fold depleted for PhastCons elements, which

250    all had more than a 1.5 fold enrichment for repeat elements (**Fig. 5A**).

251        To gain a more refined understanding of the relationship between the full-stack chromatin states

252    and conservation, we analyzed their enrichment using 100 previously defined conservation states by the

253    ConsHMM method [32]. These conservation states were defined based on the patterns of other species'

254    genomes aligning to or matching the human reference genome within a 100-way vertebrate alignment. We

255    observed 29 different conservation states maximally enriched for at least one full-stack state (**Fig. 3B,**

256    **Supp. Fig. 22-23**).

257        These states included, for example, ConsHMM state 1, a conservation state corresponding to bases

258     aligning and matching through all vertebrates and hence most associated with constraint. ConsHMM state

259     1 had >= 10 fold enrichment for exon associated full stack states TxEx1-4 and TxEnh4 (**Supp. Fig. 22**).

260     Another ConsHMM state, state 28, which is associated with moderate aligning and matching through many

261     vertebrates and strongly enriched around TSS and CpG islands, had a 44.5 and 47.8 fold enrichment for

262     TSS-associated full-stack states TSS1 and TSS2, respectively (**Supp. Fig. 22**). Additionally, this

263     conservation state is consistently the most enriched conservation state for full stack states associated with

264     flanking and bivalent promoters (**Fig. 3B, Supp. Fig. 22**). ConsHMM state 2, which has high aligning and

265     matching frequencies for most mammals and a subset of non-mammalian vertebrates and previously

266     associated with conserved enhancer regions [32], showed >2.7 fold enrichment for some full-stack enhancer

267     states for Brain (EnhWk4 and EnhA6), ESC & iPSC (EnhA17,19 and EnhWk8), neurosph (EnhWk4,

268     EnhA2,17), and mesenchymal, muscle, heart, adipose (EnhA2) (**Fig. 3B, Supp. Fig. 22**).

269        ConsHMM state 100, a conservation state associated with alignment artifacts, was 10.9 folds

270     enriched for full-stack state ZNF1, which showed 20.8 fold enrichment with ZNF genes (**Fig. 3A-B, Supp.**

271     **Fig. 22**). This is consistent with previous analysis using cell-type-specific annotations showing that

272     ConsHMM state 100 was enriched in a ZNF gene-associated chromatin state [32]. Interestingly though,

273     another full-stack state (ZNF2) that was even more strongly enriched for ZNF genes (68.6 folds), had 0.4

274     fold enrichment for ConsHMM state 100, and instead was most enriched with ConsHMM state 1 (**Fig. 3A-**

275     **B, Supp. Fig. 22**). Therefore, the full-stack annotation helped distinguish two ZNF-gene associated states

276     that are associated with distinct conservation states. As this example illustrates, the full-stack annotation

277     captured conservation state enrichments that were generally consistent with those seen in cell-type-specific

278     annotations, but could also identify additional refined enrichment patterns.

279

280     *Specific full stack states show distinct enrichments and depletions for structural variants*

281     We also analyzed the enrichment of the full-stack states for overlap with structural variants (SVs) mapped

282     in 17,795 deeply sequenced human genomes [33]. We focused on the two largest classes of SVs, deletions

283    and duplications, that were previously analyzed using 15-state cell-type-specific chromatin state models

284    [16,33]. In those analyses, enrichments were computed for 1-kb windows that were stratified based on the

285    number of cell or tissue types each state was present. ZNF gene and heterochromatin states were enriched

286    for deletions and duplications, with the enrichments being strongest when bases were annotated as those

287    states in larger numbers of cell or tissue types [33].

288         Consistent with those previous results, using the full-stack model, we observed that of the 13 states

289    that were in the top 10 maximally enriched states with either deletions or duplications (1.18 fold or greater),

290    seven were in the heterochromatin group (HET1-2,4-7,9) and one was in the ZNF gene state (ZNF2) (**Fig.**

291    **6A, Supp. Fig. 24**). The other five states included one artifact state (GapArtf2), three quiescent states

292    (Quies1-2,4) and a polycomb repressed state (ReprPC8) (**Fig. 6A**). The quiescent states Quies1-2,4, despite

293    the generally low frequencies for all marks, did have higher emission probabilities for H3K9me3 compared

294    to other chromatin marks (**Fig. 6B**). While the full-stack states showed generally consistent patterns of

295    enrichments with the analysis of [33], it allowed a more refined analysis of enrichment patterns with

296    structural variants. For example, it identified a polycomb repressed state (ReprPC8) that was enriched with

297    duplication (1.21 fold enriched) and yet depleted with deletions (5 fold depleted) (**Fig. 6A**).

298         The full-stack model was also more predictive of SV than cell-type-specific annotations. In

299    comparing with cell-type-specific annotations, the full-stack model had better AUROCs for predicting

300    locations of deletions and duplications than the 18-state model in all cases, and the 100-state cell-type-

301    specific model in all cases except for two out of 127 cell-types (**Supp. Fig. 25-26**). Additionally, we verified

302    that the full-stack model had higher AUROC in predicting duplications and deletions compared to

303    annotations obtained by ranking genomic bases based on the number of cell or tissue types that a state was

304    observed separately for each state in the 15-state model, as in the approach of [33] (**Methods**, **Supp. Fig.**

305    **27**). These results show that the full-stack annotation can uncover enrichment patterns with SVs that are

306    consistent with cell-type-specific annotations, yet highlight states with greater predictive power and offer a

307    more refined chromatin annotation of the regions enriched with SVs.

308

13

309 *Full stack-state gives insights into bases prioritized by different variant prioritization scores*

310     Various scores have been proposed to prioritize deleterious variants in non-coding regions of the

311 genome or genome-wide. These scores are based on either conservation or on integrating diverse sets of

312 genomic annotations. Though the scores all serve to prioritize variants, they can vary substantially from

313 each other and it is often not clear the differences among the types of bases that different scores prioritize.

314 To better understand the epigenomic contexts of bases that each score tends to prioritize, we analyzed the

315 full-stack state enrichment for bases they prioritize. As the scores we considered are not specific to a single

316 cell type, the full-stack states have the potential to be more informative for this analysis than cell-type-

317 specific annotations. We considered a set of 14 different variant prioritization scores that were previously

318 analyzed in the context of conservation state analysis [32]. The 14 scores for which we analyzed prioritized

319 variants in non-coding regions were CADD(v1.4), CDTS, DANN, Eigen, Eigen-PC, FATHMM-XF, FIRE,

320 fitCons, FunSeq2, GERP++, LINSIGHT, PhastCons, PhyloP, and REMM [34–46]. For each of these

321 scores, we first analyzed the full-stack state enrichments for the top 1% prioritized non-coding variants

322 relative to the background of non-coding regions on the genome (**Methods**).

323     In the top 1% prioritized non-coding bases, 19 states were among the top five most enriched states

324 ranked by at least one of the 14 scores (**Fig. 6C, Supp. Fig. 28, 29**). These 19 states include six states in

325 flanking promoter and TSS groups, three states in the bivalent promoter group, five states in enhancers and

326 transcribed enhancers groups, three states in the exon-associated transcription group, one polycomb

327 repressed state, and one DNase state (**Fig. 6C**). Seven scores (DANN, Eigen, Eigen_PC, funSeq2, CDTS,

328 CADD and REMM) had their top five enriched states exclusively associated with promoter and TSS states

329 (PromF2-5, TSS1-2, BivProm1-2,4), with enrichments ranging between 8.6 and 70 fold (**Fig. 6C**). Some

330 other scores, however, showed relatively weak enrichments or even depletions for these promoter- and

331 TSS- associated states. For example, state PromF4, which had over 30 fold enrichment for non-coding

332 variants prioritized by four different scores, had a 5-fold depletion for those prioritized by fitCons (**Fig.**

333 **6C**). Similarly, state TSS1 was in the top five most enriched states with bases prioritized by 10 scores (~ 5-

334 62 folds), including the aforementioned seven scores, yet was depleted with prioritized variants by fitCons

14

335   (~ 1.2 fold depletion) (**Fig. 6C**). Enhancer states EnhA2-3,17 were among the states in the top five most

336   enriched for FATHMM, GERP++, LINSIGHT, PhastCons, and PhyloP prioritized non-coding variants.

337   These states also showed consistent enrichments with variants prioritized by Eigen, funSeq2, CADD, and

338   REMM, though those scores showed even stronger relative enrichments for promoter states. In contrast,

339   FIRE, DANN and CDTS were depleted for prioritized variants in all these enhancer states, and Eigen_PC

340   showed both enrichments and depletions (**Fig. 6C**). FIRE and fitCons showed strong enrichment for exon

341   states (TxEx1-3), which are associated with coding regions, even though coding bases were excluded in

342   this analysis (**Fig. 6C**). FATHMM had the greatest relative enrichment for the primary DNase state

343   associated with CTCF cell type-specific chromatin states (DNase1) (~10 fold), and was the only score for

344   which this state was among the top five most enriched states (**Fig. 6C, Supp. Fig. 28**).

345        We conducted similar analyses based on top 5% and 10% prioritized non-coding variants and

346   observed relatively similar patterns of enrichments, though there did exist some differences at these

347   thresholds (**Supp. Fig. 28, 30-31**). One difference was that alignment artifact associated states GapArtf2-3

348   were among the top two states most enriched with top 5% and 10% non-coding bases prioritized by

349   FATHMM (**Supp. Fig. 28**). In addition, we analyzed top 1%, 5%, and 10% prioritized variants genome-

350   wide from the 12 scores that were defined genome-wide (**Methods**). Compared to the non-coding analysis,

351   we saw a larger number of scores that have exon-associated transcription states (TxEx1-TxEx4) among the

352   top five enriched states with top 1% variants genome-wide, while we saw no enhancer state among the top

353   five enriched states with top 1% variants by any score and only one enhancer state among the top five by

354   one score (GERP++) for top 5% and 10% variants (**Supp. Fig. 32**).

355        We verified that the full-stack annotation showed the highest AUROC in recovering the top 1%

356   non-coding variants compared to all 18-state cell-type-specific annotations for all 14 scores (**Supp. Fig.**

357   **33**). Compared to all 100-state cell-type-specific annotations, the full-stack model showed the highest

358   AUROC for 13 out of 14 scores in all 127 cell types (**Supp. Fig. 33**).

359

360   *Full-stack states show distinct enrichments and depletions for human genetic variation*

361        We next analyzed full-stack states for their enrichment with human genetic sequence variation. We

362      calculated enrichments of full-stack states with genetic variants sequenced in 15,708 genomes from

363      unrelated individuals in the GNOMAD database stratified by minor allele frequencies (MAFs) [47]. Across

364      eleven ranges of MAFs, the state enrichments ranged from a 2-fold enrichment to a 4-fold depletion (**Supp.**

365      **Fig. 34**). As expected, the state associated with assembly gaps (GapArtf1) is most depleted with variants,

366      regardless of the MAF range. At the other extreme, state Acet1, which is associated with simple repeats, is

367      the most enriched state with variants for all ten minor allele frequency (MAF) ranges that are greater than

368      0.0001, with fold enrichments between 1.8 and 2.0 (**Supp. Fig. 34**). We verified that the high enrichment

369      for state Acet1 was not specific to GNOMAD's calling of variants as it had a 2.0 fold enriched with common

370      variants from dbSNP (**Methods**) (**Supp. Fig. 34**). TSS and promoters associated states, PromF4 and TSS1-

371      2, were maximally enriched for variants in the lowest range of MAF (0 < MAF <= 0.0001), 1.5-1.7 fold.

372      The enrichment of variants for these states decreased as the MAF ranges increased, falling to 0.8-1.2 fold

373      for variants of the highest range of MAF (0.4-0.5) (**Supp. Fig. 34**). The high enrichment for states PromF4

374      and TSS1-2 for rare variants is consistent with these states having high enrichment for CpG islands (75-

375      100 fold) (**Fig. 3A**) and the high mutation rate for CG dinucleotides [48]. We observed a similar though

376      weaker pattern of decreasing enrichments for increasing MAF for other states associated with

377      transcriptional activities, enhancers, DNase, or promoters (**Supp. Fig. 26)**. This pattern was not observed

378      in most states from other groups such as heterochromatin, polycomb repressed, quiescent, and acetylations

379      only (**Supp. Fig. 26**).

380        To better identify states with a depletion of common variants that are more likely due to selection,

381      we ranked states based on their ratios of enrichments for the rarest variants (MAF < 0.0001) relative to the

382      most common variants (MAF 0.4-0.5) (**Fig. 6D**). The states with the highest ratio included a number of

383      flanking promoter (PromF3-4) and exon-transcription states (TxEx1,2,4) that were also associated with

384      strong sequence conservation across species (**Fig. 6D, Fig. 3B**). These results are consistent with previous

385      analyses supporting a depletion of common human genetic variation in evolutionary conserved regions

386      [49]. States associated with assembly gaps and alignment artifacts (GapArtf1-3), quiescent (Quies3), or

387    acetylations and simple repeats (Acet1) were most depleted for rare variants relative to the common variant

388    enrichment (**Fig. 6D**).

389

390    *Full-stack states show enrichment for phenotype-associated genetic variants*

391         We next analyzed the relationship between the full-stack states and phenotypic associated genetic

392    variants. We first evaluated the enrichment of the full-stack state for variants curated into the Genome-wide

393    Association Study (GWAS) catalog relative to a background of common variation [50] (**Methods**). This

394    revealed six states with at least a two-fold enrichment (**Supp. Fig. 35**). Four of these states, TxEx1-2,4 and

395    TxEnh4, were all transcription associated states that are $\geq$ 10-fold enriched with coding sequences and

396    $\geq$=11 fold for ConsHMM state 1, associated with the most constraint in a sequence alignment of 100

397    vertebrates (**Fig. 3B**). This observation is consistent with previous results that GWAS catalog variants show

398    enrichments for coding sequence and sequence constrained bases [32,49,51]. The other two states with

399    greater than two-fold enrichment for GWAS catalog variants relative to common variants were two

400    promoter states, PromF2-3 (**Supp. Fig. 35**). On the other hand, four states were more than two-fold depleted

401    for GWAS catalog variants, and were associated with artifacts (GapArtf2-3), or quiescent and polycomb

402    repressed states with weak signals of H3K9me3 (Quies5) or H3K27me3 (ReprPC8) (**Supp. Fig. 35**).

403         We also analyzed the full-stack state enrichments for fine-mapped variants previously generated

404    from a large collection of GWAS studies from the UK Biobank database and other public databases [52].

405    Specifically, we considered separately the fine mapped variants from two fine-mapping methods, CAVIAR

406    [53] and FINEMAP [54], for 3052 traits. For each method and trait, we identified the single variants that

407    had the greatest probability of being causal at a set of distinct loci, and computed the enrichment of these

408    variants for the full-stack states relative to a background of common variants (**Methods**).

409         Fold enrichment results of full-stack states for the most likely causal variants were highly consistent

410    between fine-mapping methods (FINEMAP and CAVIAR) (**Supp. Fig. 36**). The ten states maximally

411    enriched with fine-mapped variants relative to common variants, which were the same states by both

412    methods, included five states associated with flanking and bivalent promoter activities (PromF2-5,

17

413    BivProm4), an enhancer state in blood and thymus (EnhA9) and an enhancer state in most other cell types

414    except blood (EnhA1), and three highly conserved transcription-associated states (TxEnh4,6, TxEx4) (**Fig.**

415    **6E**). Notably, five of 10 states maximally enriched with fine-mapped variants, PromF2-5, BivProm4, were

416    associated with promoter regions and also among the 19 states most enriched with top 1% prioritized

417    variants by at least two of the 14 different variant prioritization scores (**Fig. 6E, C**). These results show that

418    there are agreements in the types of chromatin states preferentially overlapped by phenotype-associated

419    fine mapped variants and variants predicted to have greater effects based on variant prioritization

420    scores.  We also confirmed that the full-stack model consistently resulted in higher AUROC in predicting

421    locations of fine-mapped variants within a background of common variants, compared to the 18-state and

422    100-state cell-type-specific annotations in all cell types (**Supp. Fig. 37-38**).

423

424    *Full-stack states show enrichments for cancer-associated variants*

425         In addition to investigating germline variants, we also investigated the enrichment of full-stack

426    states for somatic variants identified from whole genome sequencing of cancer samples. We analyzed data

427    of variants from four cancer types that have the largest number of somatic variants in the COSMIC database

428    [55]: liver, breast, pancreas and haematopoietic_and_lymphoid_tissue (**Methods**). Sixteen states were

429    among the top 10 most enriched with at least one type of cancer's associated variants (1.2-1.4 fold in breast

430    cancer, 1.2-5.6 fold in lymphoid cancer, 1.2-5.4 in liver cancer, 1.4-4.2 in pancreas cancer) (**Fig. 6F**).

431    Among these states, 15 states showed higher signals of H3K9me3 compared to most other chromatin marks,

432    including seven states in heterochromatin group (HET1-2, 4-7,9), four states in quiescent group with weak

433    emissions of H3K9me3 (Ques 1-2,4-5), one state in the polycomb repressed group with weak signals of

434    H3K9me3 and H3K27me3 (ReprPC8), one state in the acetylation group with signals of H3K9me3 and

435    various acetylation marks (Acet1), two artifact-associated states with higher signals of H3K9me3 and

436    DNase relative to other marks (GapArtf2-3) (**Fig. 6G**). These results are consistent with previous findings

437    on an association of H3K9me3 and somatic cancer-associated variants [56,57].  Acet1 was also the state

438    most enriched with simple repeats, dbSNP 151 common variants, and variants of ten ranges of MAF from

18

439      GNOMAD (**Fig. 5C, Supp. Fig. 34**). Notably, the GapArtf2-3 states, associated with satellite repeat

440      enrichments (**Fig. 5C, Supp. Fig. 20**), were the top two most enriched states with somatic variants

441      associated with liver, pancreas and haematopoietic and lymphoid tissue cancers with 2.0-5.6 enrichment

442      fold (**Fig. 6F, Supp. Fig. 39**). We note that the association between the full-stack annotations and presence

443      of cancer variants is stronger than for the 18-state and 100-state cell-type-specific chromatin state

444      annotations for all four cancer types, as evidenced by the higher AUROC of the full-stack annotation at

445      predicting somatic variants (**Supp. Fig. 40-41**).

446

447      **Discussion**

448      We demonstrated a large-scale application of the stacked modeling approach of ChromHMM using

449      over a thousand epigenomic datasets to annotate the human genome. In the datasets, 32 chromatin marks

450      and 127 reference epigenomes were represented. We note that even though not every chromatin mark was

451      profiled in every reference epigenome we were still able to directly apply the stacked modeling to such

452      data. We conducted extensive enrichment analyses of the states with various other genomic annotations and

453      datasets including gene features, genetic variation, repetitive elements, comparative genomic annotations,

454      and bases prioritized by different variant prioritization scores. These analyses highlighted diverse

455      enrichment patterns of the states. Using these enrichments along with the model parameters, we provided

456      a detailed characterization of each of the 100 states in the model.

457      We grouped these 100 states into 16 groups that included promoters, enhancers, transcribed

458      regions, polycomb repressed regions, zinc finger genes among others. We also highlighted important

459      distinctions among states within the groups. In many cases, identifying these distinctions was enabled by

460      the full-stack modeling using data from multiple cell types for genome annotation. For example, we

461      identified enhancer and repressive states that were active in different subsets of cell types. We also

462      highlighted how different states in some of the groups such as those associated with transcribed and ZNF

463      genes showed distinct enrichments for conservation states. Overall, the full-stack model showed enrichment

464      patterns supporting observations held for cell-type-specific annotations, yet it provided more detailed
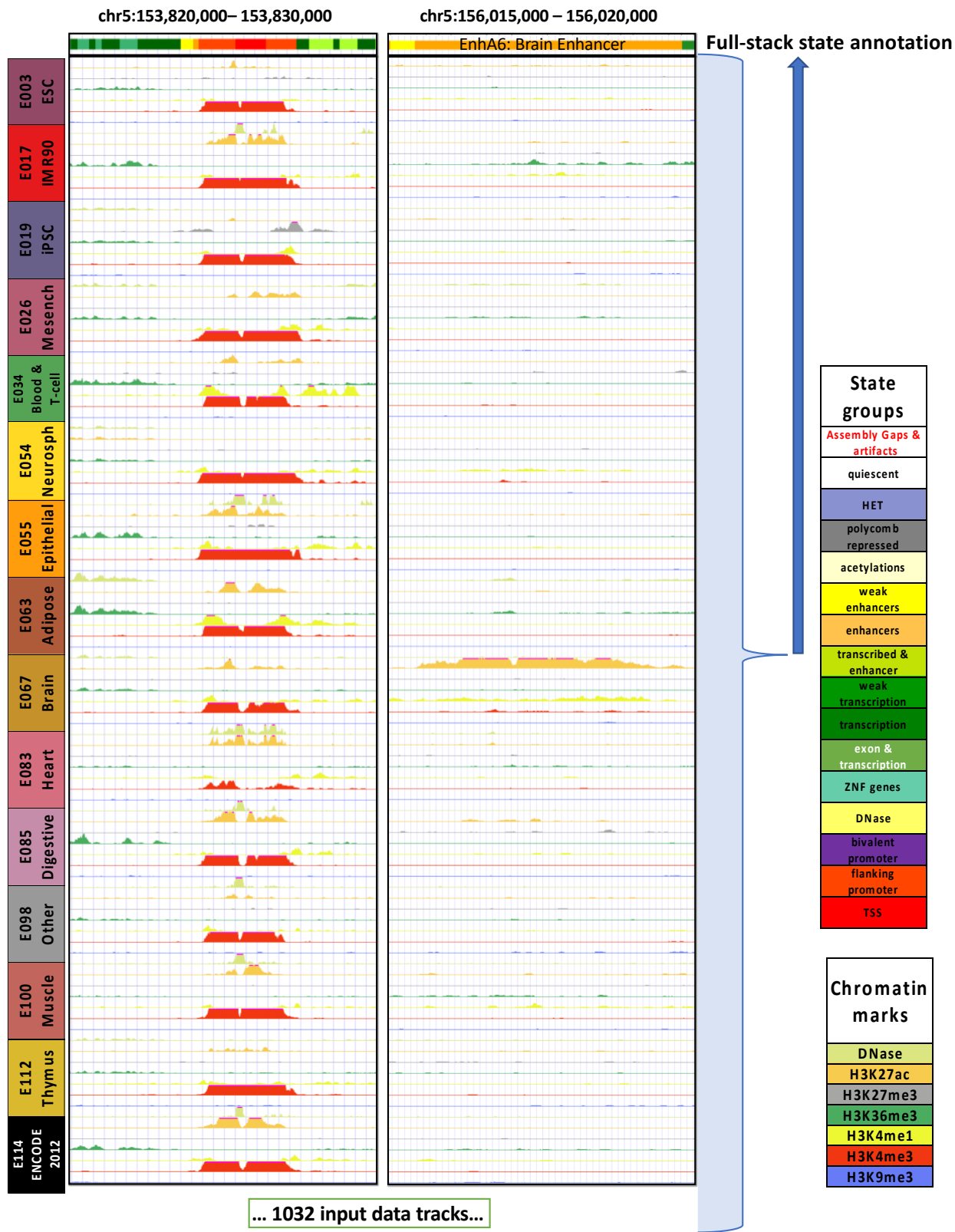
19

465    stratification of genomic regions into chromatin states with heterogeneous associations with other genomic

466    information.

467         The full-stack modeling has advantages to commonly used cell-type-specific chromatin state

468    annotations in several respects. First, the full-stack model is not specific to one cell or tissue type and thus

469    is able to provide a unified view of all the data and directly uncover states that correspond to constitutive

470    or cell-type-specific activities. Second, the full-stack annotation consistently showed better recovery of

471    various genomic features compared to cell-type-specific annotations. This improvement is expected since

472    full-stack models can leverage information from multiple cell types for genome annotations. Third, in cases

473    where it is not desirable to focus on only one specific cell or tissue for analysis, the full-stack modeling can

474    bypass the need to pick one such cell or tissue type or to consider a large number of different cell-type-

475    specific chromatin state annotations simultaneously. Such cases may arise when studying other genomic

476    information that is not inherently cell-type-specific such as genome variation and sequence conservation.

477         However, we emphasize that the stacked modeling approach should be considered a complement

478    to and not a replacement of the cell-type-specific annotations, which have their own advantages. Cell-type-

479    specific annotations may be preferable when one is interested in a specific cell type or in directly comparing

480    the chromatin state maps among individual cell types. Additionally, the cell-type-specific chromatin states

481    have fewer parameters and thus can be easier to interpret relative to stacked model states.

482         We expect many applications of the full-stack annotations that we generated here. The full-stack

483    annotation can be used as a resource to interpret genetic variation. A possible avenue for future work is to

484    incorporate the full-stack annotation into scoring methods to better predict genetic variants' phenotypic

485    influences. Future work could apply the stacked modeling approach to even larger sets of data that are

486    accumulating in human as well as large datasets in key model organisms such as mouse. This work provides

487    a new annotation resource for studying the human genome, non-coding genetic variants, and their

488    association with diseases.

489

490

491 **Figure 1: Illustration of full-stack modeling annotations.** The figure illustrates the full-stack modeling

492 at two loci. The top track shows chromatin state annotations from the full-stack modeling colored based on

493 the legend at right. Below it are signal tracks for a subset of the 1032 input datasets. Data from seven

494 (DNase I hypersensitivity, H3K27me3, H3K36me3, H3K4me1, H3K4me2, H3K4me3, and H3K9me3) of

495 the 32 chromatin marks are shown, colored based on the legend at right. These data are from 15 of the 127

496 reference epigenomes each representing different cell and tissue groups. The loci on left highlights a

497 genomic region for which a portion is annotated as constitutive promoter states (TSS1-2). The loci on right

498 panel highlights a region for which a portion is annotated as a brain enhancer state (EnhA6), which has high

499 signals of H3K27ac in reference epigenomes of the group Brain.

500

501

502

503

504

505

506

507

508

509

510

511

512

513

514

515

516

543 **Figure 2: Full-stack state emission parameters. (A)** Each of the 100 rows in the heatmap corresponds to

544 a full-stack state. Each of the 1032 columns corresponds to one experiment. For each state and each

545 experiment, the heatmap gives the probability within the state of observing a binary present call for the

546 experiment's signal. Above the heatmap there are two rows, one indicating the cell or tissue type of the

547 experiment and the other indicating the chromatin mark. The corresponding color legends are shown

548 towards the bottom. The states are displayed in 16 groups with white space between each group. The states

549 were grouped based on biological interpretations indicated by the color legend at the bottom. Full

550 characterization of states is available in **Supplementary Data**. The model's transition parameters between

551 states can be found in **Supp. Fig. 6.** Columns are ordered such that experiments profiling the same

552 chromatin marks are next to each other.

553 **(B)** Each row corresponds to a full-stack state as ordered in (A). The columns correspond to the top 10

554 experiments with the highest emission value for each state, in order of decreasing ranks, colored by their

555 associated chromatin marks as in (A).

556 **(C)** Similar to **(B)**, but experiments are colored by the associated cell or tissue type group. We noted on the

557 right the cell or tissue groups of some cell-type-specific enhancer states.

558

559

560

561

562

563

564

565

566

567

568

Figure panels:
- **A** Enrichments with external annotations
- **B** Associated consHMM state
- **C** Avg. gene exp
- **D** Annotated TES neighborhood enrichments
- **E** Annotated TSS neighborhood enrichments
- **F** Enrichments with open CTCF elements

ConsHMM states legend:
- High align and match frequencies for a few primates
- High align and march for mammals, but missing notable subsets
- High align and match frequencies for primates
- Putative artifact
- High align and match frequency for all vertebrates
- High align and match for mammals
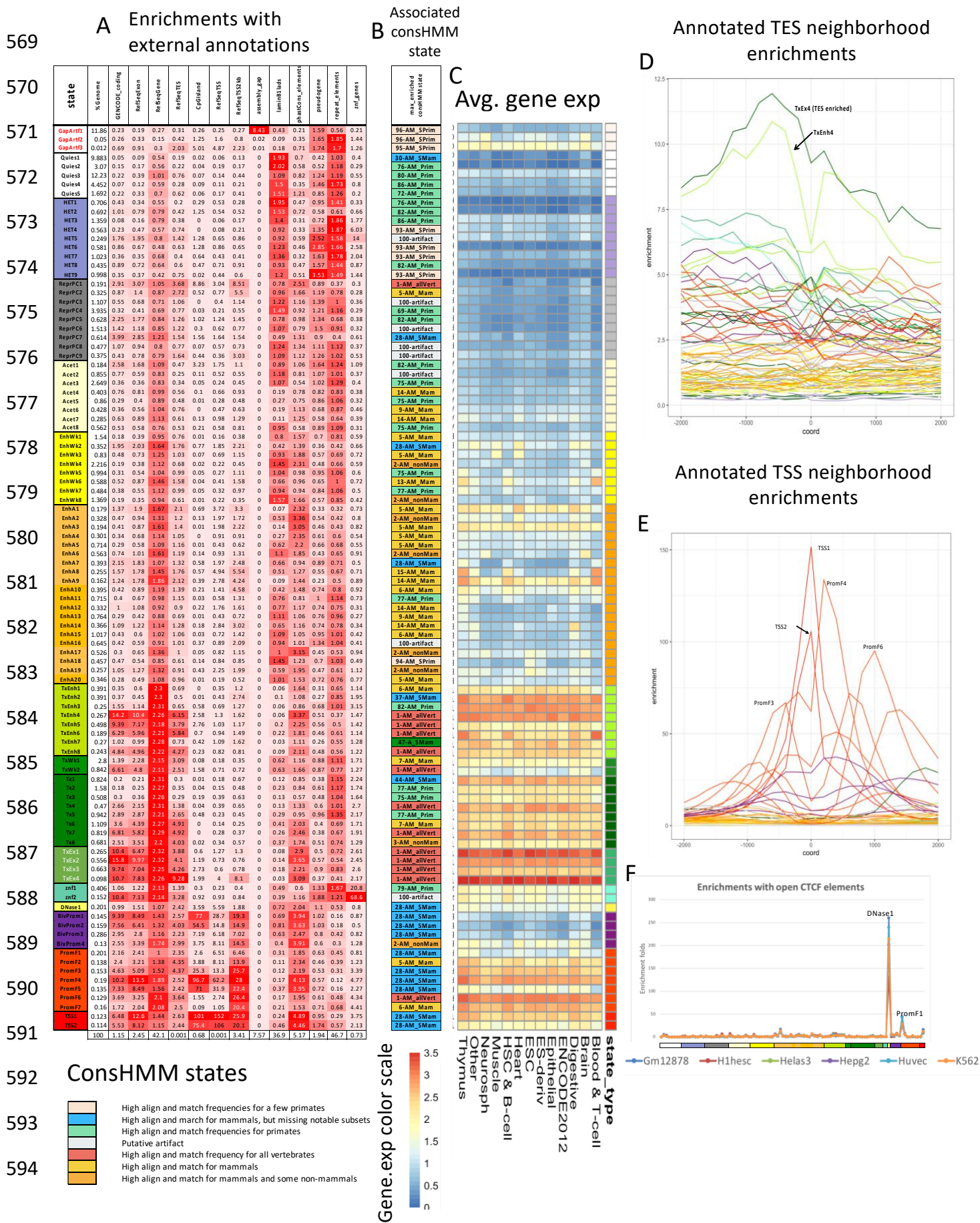- High align and match for mammals and some non-mammals

595    **Figure 3: Full stack states enrichments for external genomic annotations. (A)** Fold enrichments of full-

596    stack states with external genome annotations (**Methods**). Each row corresponds to a state and each column

597    corresponds to one external genomic annotation: CpG Islands, Exons, coding sequences, gene bodies,

598    transcription end sites (TES), transcription start sites (TSS), TSS and 2kb surrounding regions, lamina

599    associated domains (laminB1lads), assembly gaps, annotated ZNF genes, repeat elements and PhastCons

600    constrained element. The color is normalized to range from minimum values (white) to maximum values
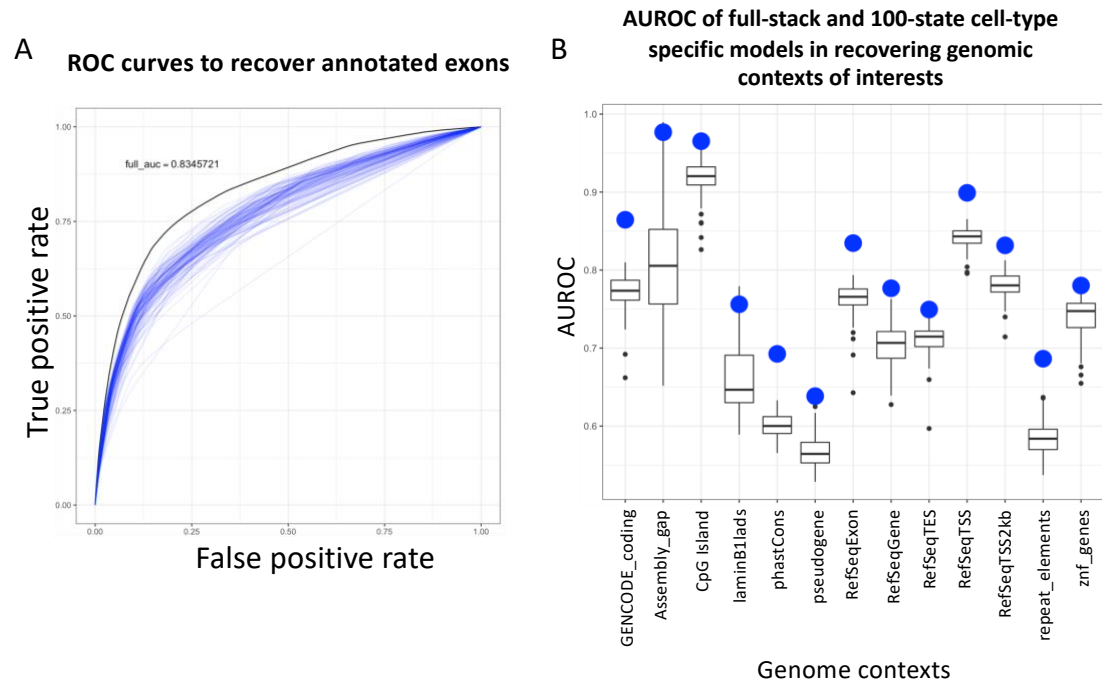
601    (red) within each column.

602    **(B)** Each row indicates the ConsHMM state [32] that has highest enrichment fold in each full-stack state as

603    ordered in **(A).** Legends of the ConsHMM state groups indicated with different colors are shown below the

604    heatmap in **(A).**

605    **(C)** Average weighted expression of genes that overlap each full-stack state in different groups of cells

606    (**Methods**). Each column corresponds to a cell group indicated at the bottom. Each row corresponds to a

607    state, as ordered in **(A).**

608    **(D-E)** Positional enrichments of full-stack states relative to annotated **(D)** transcription end sites (TES) and

609    **(E)** transcription start sites (TSS). Positive coordinate values represent the number of bases downstream in

610    the 5' to 3' direction of transcription, while negative values represent the number of bases upstream. Each

611    line shows the positional enrichments in a state. Lines are colored as indicated in **(A).**

612    **(F)** Enrichments of full-stacks states with cell-type-specific chromatin states associated with CTCF and

613    open chromatin, but limited histone modifications in six cell types [30] (**Methods**). The six cell types are

614    indicated along the bottom of the figure. States are displayed horizontally in the same order as **(A)**. The

615    DNase1 state showed the strongest enrichment for the cell-type-specific chromatin states associated with

616    CTCF and open chromatin in all six cell types.

617

26

**A** ROC curves to recover annotated exons

**B** AUROC of full-stack and 100-state cell-type specific models in recovering genomic contexts of interests

618

**Figure 4: Full-stack annotation's recovery of external genome annotations**

**(A)** ROC curves in recovering annotated exons by full-stack annotation (solid black line) and 127 cell-type-specific 100-state annotations (blurred blue lines). Full stack model yielded the highest AUROC (0.83).

**(B)** A comparison of the AUROC for full-stack annotation and cell-type-specific models for recovering positions of different external annotations. Each box plots show the range of AUROC of 100-state cell-type-specific chromatin state annotations for recovery of one external annotation and the large blue point shows the AUROC for the full-stack annotation. The external annotations in order were coding sequences, assembly gaps, CpG Islands, lamina associated domains, phastCons conserved elements, pseudogenes, exons, gene bodies, transcription end sites (TES), transcription start sites (TSS), TSS and 2kb surrounding regions, repeat elements, annotated ZNF genes. These annotations are similar to **Fig. 3A**. ROC curves corresponding to these AUROC values can be found in **Supp**. **Fig. 18**.

27

**Figure 5: Full-stack states enrichments with conserved elements and repeat classes.**

**(A)** The first ten rows show the states most enriched with PhastCons elements and concurrently least enriched with RepeatMasker repeat elements, ordered by decreasing enrichments with PhastCons elements. The bottom ten rows show the states most enriched with repeat elements and concurrently least enriched with PhastCons elements, ordered by increasing enrichments with repeat elements. The columns from left to right list the state ID, the percent of the genome that each state covers, and the fold enrichments for repeat elements and PhastCons elements.

**(B)** Heatmap of the state emission parameters from **Fig. 2A** for the subset of states highlighted in panel (A). The colors are the same in **Fig. 2A**.

**(C)** Fold enrichments of full-stack states with different repeat classes (**Methods**). Rows correspond to states and columns to different repeat classes. Only states that are most enriched with at least one repeat class are shown. Fold enrichment values that are maximal for a given are shown in dark red. Other fold enrichments greater than one are shaded light red.

28

A

| state | % genome | deletion | duplication |
|---|---|---|---|
| GapArtf2 | 0.05 | 1.27 | 1.55 |
| Quies1 | 10.7 | 1.62 | 1.21 |
| Quies2 | 3.33 | 1.56 | 1.25 |
| Quies4 | 4.83 | 1.22 | 1.13 |
| HET1 | 0.77 | 1.57 | 1.3 |
| HET2 | 0.75 | 0.99 | 1.45 |
| HET4 | 0.61 | 1.52 | 1.45 |
| HET5 | 0.27 | 1.22 | 1.2 |
| HET6 | 0.63 | 1.34 | 1.36 |
| HET7 | 1.11 | 1.18 | 1.12 |
| HET9 | 1.08 | 1.33 | 1.44 |
| ReprPC8 | 0.52 | 0.16 | 1.21 |
| znf2 | 0.16 | 1.09 | 1.25 |

B



C

| state | % genome | FIRE | fitCons | FATHMM | GERP | LINSIGHT | PhastCons | phyloP | DANN | CDTS | CADD | REMM | Eigen | Eigen_PC | funSeq2 | long state annotations |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| ReprPC1 | 0.2 | 0 | 1.1 | 4.4 | 3.9 | 3.9 | 3.6 | 3.8 | 0.7 | 11 | 4.9 | 5 | 5.2 | 4.3 | 1.5 | ReprPC (except: PromBiv in ESC); H3K27me3 strong and H3K4me1 weak |
| EnhA2 | 0.36 | 0.3 | 1.9 | 4.5 | 5.6 | 5.2 | 4.8 | 4.2 | 0.2 | 0.8 | 6.3 | 7.4 | 6.2 | 0.8 | 4.6 | enhancers in mesenchymal, muscle, heart, neurosph, adipose. DNase in ESC and iPSC |
| EnhA3 | 0.21 | 0.3 | 9.2 | 3.2 | 4.5 | 4.5 | 3.9 | 3.4 | 0.3 | 0.8 | 5.6 | 5.6 | 8 | 7.9 | 12 | strong enhancers in most cells, weaker in blood and ESC&iPSC |
| EnhA17 | 0.58 | 0.2 | 1 | 4.8 | 5.7 | 4.8 | 4.7 | 4.3 | 0.3 | 0.6 | 5.9 | 5.8 | 5.5 | 0.1 | 2.3 | ESC, iPSC, Neurosph enhancers and some differentiated |
| TxEnh4 | 0.25 | 15 | 17 | 2.9 | 2 | 2.2 | 1.5 | 1.8 | 0.7 | 2.8 | 2.1 | 2.8 | 2 | 4.1 | 7.8 | H3K36me3 strong and some enhancers (H3K27ac, H3K4me1); exons; near TES; 11_TxEnh3p all cell types |
| TxEnh8 | 0.25 | 4.1 | 9 | 2 | 2.2 | 2.3 | 1.6 | 1.8 | 0.4 | 0.8 | 1.9 | 2.6 | 2.3 | 2.5 | 7.7 | Transcribed enhancers 3'; enhancers Myostat, IMR90, Mesench, Epithelial |
| TxEx1 | 0.26 | 9.9 | 10 | 2.5 | 2.1 | 1.9 | 1.7 | 2 | 0.6 | 1.2 | 2 | 2.2 | 1.8 | 0.1 | 3.2 | H3K79me2 and H3K36me3; exon; 6_Tx |
| TxEx2 | 0.5 | 13 | 14 | 3.2 | 1.9 | 1.8 | 1.4 | 1.8 | 0.8 | 2.5 | 1.6 | 1.9 | 1.1 | 0.3 | 2.1 | H3K36me3 strongest; exons; Tx3p |
| TxEx3 | 0.65 | 13 | 8.7 | 1.8 | 1 | 1.1 | 0.9 | 1.2 | 0.7 | 2.1 | 0.9 | 1.1 | 0.7 | 0.8 | 1.1 | H3K36me3 strong; exon; Tx3p |
| DNase1 | 0.22 | 0.3 | 3 | 10 | 2.2 | 3 | 2.4 | 2.5 | 2.2 | 2.3 | 2.8 | 3.2 | 4.3 | 10 | 7.1 | DNase I only; CTCF, Candidate Insulator |
| BivProm1 | 0.14 | 0.2 | 1.1 | 4.5 | 2.8 | 4.1 | 4.1 | 3.8 | 8.8 | 52 | 8.7 | 13 | 11 | 32 | 14 | bivalent promoter- more balanced H3K4me3/ H3K27me3 |
| BivProm2 | 0.16 | 0 | 1.3 | 4.8 | 3.5 | 4.2 | 4.2 | 4.1 | 4.8 | 41 | 7.3 | 8.9 | 8.6 | 16 | 6.9 | bivalent promoter- stronger on H3K27me3 |
| BivProm4 | 0.14 | 0.5 | 1.7 | 5.4 | 6.6 | 5.8 | 5.7 | 5.4 | 0.7 | 7.8 | 8.4 | 9.7 | 9.8 | 6.5 | 8.6 | 24_ReprPC in ENCODE, Blood & Tcell, Digestive, HSC &B-cell, Sm.Muscle. 23_PromBiv flanking in other cell types. H3K27me3 and H3K4me1 |
| PromF2 | 0.15 | 4.1 | 5.5 | 2.6 | 2.6 | 3.5 | 2.7 | 2.5 | 0.9 | 11 | 5 | 5.3 | 10 | 40 | 21 | H3K4me1, H3K4me2, H3K4me3, DNase, acetylations promoter flank upstream bias |
| PromF3 | 0.16 | 11 | 0.8 | 1.2 | 1.4 | 2.6 | 2 | 1.8 | 2.3 | 35 | 4.2 | 6.7 | 16 | 64 | 33 | H3K4me2, H3K4me3, H3K4me1(weaker than me3), DNase, acetylations - flanking tss upstream and downstream |
| PromF4 | 0.19 | 13 | 0.2 | 2.4 | 2.6 | 4.5 | 4 | 3.4 | 9.8 | 56 | 8.7 | 19 | 32 | 73 | 38 | H3K4me2, H3K4me3 limited H3K4me1, heavily acetylated - flanking tss downstream bias |
| PromF5 | 0.14 | 0.8 | 1.1 | 3.7 | 3 | 4.2 | 4.4 | 3.9 | 8.9 | 48 | 9.1 | 17 | 16 | 52 | 21 | flanking promoter; 2_PromU in most cell types; stronger on H3K4me3 |
| TSS1 | 0.13 | 6 | 0.8 | 3 | 3.4 | 5.3 | 6.1 | 4.8 | 17 | 41 | 12 | 24 | 26 | 61 | 30 | TSS more acetylated and active; 1_TssA in all cell types |
| TSS2 | 0.12 | 1.5 | 2 | 4.2 | 2.5 | 3.7 | 5.3 | 4 | 19 | 23 | 9.6 | 14 | 16 | 40 | 16 | TSS (22_PromP in iPSC) |

D



Log of enrichment ratios with least MAF variants vs. highest MAF variants

Full-stack states

| state | enrichment ratio |
|---|---|
| GapArtf2 | 0.58 |
| GapArtf3 | 0.65 |
| Acet1 | 0.75 |
| Quies3 | 0.83 |
| GapArtf1 | 0.85 |
| TxEx1 | 1.73 |
| PromF3 | 1.78 |
| TxEx4 | 1.78 |
| TxEx2 | 1.83 |
| PromF4 | 1.92 |

E

| state | % background | caviar_lead_snps | finemap_lead_snps |
|---|---|---|---|
| EnhA1 | 0.17 | 2.92 | 2.93 |
| EnhA9 | 0.15 | 2.72 | 2.72 |
| TxEnh4 | 0.23 | 2.93 | 2.92 |
| TxEnh6 | 0.17 | 2.65 | 2.65 |
| TxEx4 | 0.08 | 3.36 | 3.36 |
| BivProm4 | 0.12 | 2.86 | 2.86 |
| PromF2 | 0.12 | 3.05 | 3.04 |
| PromF3 | 0.13 | 3 | 3.01 |
| PromF4 | 0.19 | 3.08 | 3.08 |
| PromF5 | 0.13 | 2.88 | 2.89 |

rank of most enriched states: 1, 2, 3, 4, 5, >5

rank of least enriched: 1, 2, 3, 4, 5

F

| state | % genome | breast | haematopoietic and lymphoid tissue | liver | pancreas |
|---|---|---|---|---|---|
| GapArtf2 | 0.05 | 0.76 | 4.88 | 2.07 | 4.09 |
| GapArtf3 | 0.01 | 1.35 | 5.58 | 5.38 | 4.22 |
| Quies1 | 10.84 | 1.21 | 1.69 | 1.54 | 1.57 |
| Quies2 | 3.36 | 1.26 | 1.41 | 1.49 | 1.75 |
| Quies3 | 13.37 | 0.97 | 1.23 | 0.95 | 0.94 |
| Quies4 | 4.86 | 1.18 | 1.25 | 1.09 | 1.23 |
| Quies5 | 1.85 | 1.34 | 0.96 | 0.79 | 0.98 |
| HET1 | 0.77 | 1.25 | 1.38 | 1.53 | 1.95 |
| HET2 | 0.75 | 1.29 | 0.92 | 1.20 | 1.68 |
| HET4 | 0.61 | 0.87 | 1.14 | 0.89 | 1.47 |
| HET5 | 0.27 | 1.04 | 0.97 | 1.23 | 1.32 |
| HET6 | 0.63 | 1.37 | 1.29 | 1.31 | 1.89 |
| HET7 | 1.12 | 1.19 | 1.17 | 1.03 | 1.37 |
| HET9 | 1.08 | 1.44 | 1.34 | 1.18 | 1.56 |
| ReprPC8 | 0.52 | 1.24 | 0.75 | 0.64 | 0.77 |
| Acet1 | 0.20 | 0.83 | 2.67 | 1.23 | 1.37 |

G



644

645 **Figure 6: Full-stack states' relationship with human genetic variants.**

646 **(A)** Enrichments of full-stack states with duplications and deletions from [33]. Only states that are in the

647 top ten most enriched states are shown. Top five fold-enrichments for each class of structural variants are

648 colored in increasing darker shades of red for higher ranked enrichments. Enrichment values below one,

649 corresponding to depletions, are colored yellow. The columns from left to right are the state label, percent

650 of genome the state covers, the fold enrichment for deletions, and fold enrichment for duplications.

651 **(B)** Emission probabilities corresponding to states in **(A).** The coloring is the same as **Fig. 2A**. The figure

652 highlights how states most associated with structural variants generally had higher emission of H3K9me3

653 compared to other chromatin marks.

654 **(C)** Enrichments of full-stack states with top 1% prioritized bases in the non-coding genome by 14 variant

655 prioritization scores previously analyzed [32]. Only states that are among the top five most enriched states

656 by at least one score are shown. The top five enrichment values for each score are colored in increasing

657 darker shades of red for higher ranked enrichment values. Enrichment values below one, corresponding to

658 depletions, are colored in yellow. The columns from left to right are the state label, percent of the genome

659 covered, the 14 score enrichments, and a detailed description of the state.

660 **(D)** Log base 10 of ratios of states' enrichment with GNOMAD variants with the lowest MAFs (< 0.0001)

661 vs. GNOMAD variants with the highest MAFs (0.4-0.5). States are ordered as in **Fig. 2A**. Top five states

662 that with the highest and lowest enrichment ratios are labeled to the right.

663 **(E)** States most enriched with fine-mapped phenotypic variants against the background of common variants.

664 Fine-mapped phenotypic variants were identified by either CAVIAR [53] or FINEMAP [54] (**Methods**).

665 **(F)** State enrichments with somatic mutations associated with four cancer types in the non-coding genome.

666 Only states that are among the ten most enriched with variants from at least one cancer type are shown.

667 States in the top five are colored according to their ranks. The top five enrichment values for each cancer

668 type are colored in increasing darker shades of red for higher ranked enrichment values. The columns are

669 the state label, the percent of the genome the state covers, and the fold enrichments of variants from breast,

670 haematopietic and lymphoid, liver, and pancreas cancer types.

30

671    **(G)** Emission probabilities corresponding to states in (G), as subsetted from **Fig. 2A**. The coloring is the

672    same as **Fig. 2A**. The figure highlights how states with the greatest enrichments for cancer-associated

673    variants tend to have higher emission probabilities for H3K9me3 compared to other chromatin marks.

674

675

676

677

678

679

680

681

682

683

684

685

686

687

688

689

690

691

692

693

694

695    **Methods**

696    **Input data and processing**

697          We obtained coordinates of reads aligned to Human hg19 in .tagAlign format for the consolidated

698    epigenomes as processed by the Roadmap Epigenomics Consortium from

699    https://egg2.wustl.edu/roadmap/data/byFileType/alignments/consolidated/. In total we obtained data for

700    1032 experiments and their corresponding input control data. The experiments correspond to 127 reference

701    epigenomes, 111 of which were generated by the Roadmap Epigenomics Consortium and 16 were generated

702    by the ENCODE Consortium. Of the 1032 experiments, 979 were of ChIP-seq data targeting 31 different

703    epigenetic marks and 53 were of DNase-seq (**Sup Fig. 2**). For each of the 127 reference epigenomes there

704    was a single ChIP-seq input control experiment. For the 53 reference epigenomes that had a DNase-seq

705    experiment available there was an additional DNase control file.

706          We next binarized the data at 200 base pair resolution using the BinarizeBed command of

707    ChromHMM (v.1.18). To apply BinarizeBed in stacked mode we generated a cell_mark_file input table for

708    ChromHMM with four tab-delimited columns. The first column had the word 'genome' for all datasets, the

709    second column contained entries of the form '<EID>-<mark>' where 'EID' is the epigenome ID and 'mark'

710    is the mark name, the third column specifies the name of the corresponding file with aligned reads, and the

711    fourth column is the name of the file with the corresponding control reads. Each row in the table corresponds

712    to one of the 1032 experiments.

713

714          In order to reduce the memory and time needed to execute BinarizeBed on a large number of

715    datasets, we split the cell_mark_file table into 104 smaller tables with each table having at most 10 entries

716    corresponding to at most 10 datasets to be processed. This was done with a custom script, but the same

717    functionality has been included with the '-splitcols' and '-k' flags of BinarizedBed in ChromHMM

718    v1.22.  We then ran BinarizeBed in parallel for each of these smaller cell_mark_file tables and generated

719    output into separate sub-directories. We ran BinarizeBed with the option '-gzip' which generates gzipped

720    files.

32

721      To merge data from the 104 subdirectories from the previous step into files containing binarized

722    data of all experiments, we ran the command 'MergeBinary', which we added in v1.18 of ChromHMM.

723    We ran the command with the options '-gzip -splitrows'. The '-splitrows' option generates multiple files

724    of merged binarized data for each chromosome, where, under the default settings that we used, each file

725    contains data for a genomic region of at most 1MB. Splitting each chromosome into smaller regions allows

726    the model learning step of ChromHMM to scale in terms of memory and time to the large number of input

727    data tracks (i.e. features) that we were using. We used chr1-22, chrX, chrY, and chrM in the binarization

728    and model learning.

729

730    **Training full-stack model and generating genome-wide state annotations**

731      We learned the full-stack chromatin state model for the 1032 datasets using the LearnModel

732    command of ChromHMM (v1.18). This version of ChromHMM includes several options that we added to

733    improve the scalability when training with large numbers of features. One of these features was to randomly

734    sample different segments of the genome for training during each iteration, instead of training on the full

735    genome. This sampling strategy was previously used by ConsHMM [32], which was built on top of

736    ChromHMM. To learn the full-stack model with input data processed as outlined above, we used

737    ChromHMM's LearnModel command with the options '-splitrows -holdcolumnorder -pseudo -many -p 6

738    -n 300 -d -1 -lowmem -gzip'.

739      The '-splitrows' flag informs ChromHMM that binarized data for a chromosome is split into

740    multiple files, which reduces the memory requirements and allows ChromHMM to select a subset of the

741    genome to train on for each iteration. The '-holdcolumnorder' flag prevents ChromHMM from reordering

742    the columns of the output emission matrix, which saves time when there is a large number of features.

743      The '-pseudo' flag specifies that in each update of model parameters, ChromHMM adds a pseudo

744    count of one to the numbers of observations of transition between each pair of states, presence and absence

745    of each mark from each state, and initial state assignments of the training chromatin state sequence. This

746     prevents model parameters from being set to zero, which is needed for numerical stability when some

747     features are sparse and ChromHMM does not train on the full genome in each iteration.

748          The '-many' flag specifies ChromHMM to use an alternative procedure for calculating the state

749     posterior probabilities that is more numerically stable when there are a large number of features. The

750     procedure is designed to prevent all states from having zero posterior probability at any genomic position,

751     which can happen due to the limits of floating-point precision. The procedure does this by leveraging the

752     observation that only the relative product of emission probabilities across states are needed at each position

753     to determine the posterior probabilities. Specifically, for each position, the procedure initializes the product

754     of emission probabilities for all features, i.e. the emission product, from each state to one. For each feature,

755     the procedure then multiplies the current emission products from each state by the emission probability of

756     the feature in the state, and divides all the resulting products by their maximum to obtain updated emission

757     products. We iteratively repeat these steps of multiplication and normalization until all features have been

758     included into the calculation of relative emission products across states.

759          The '-p 6' flag specifies to ChromHMM to train the model in parallel using 6 processors. The '-n

760     300' flag specifies to ChromHMM to randomly pick 300 files of binarized data, corresponding to 300

761     regions of 1 MB (or less if the last segment of the chromosome was selected) for training in each iteration.

762     The '-d -1' option has ChromHMM not require an evaluated likelihood improvement between iterations to

763     continue training since likelihood decreases are expected as on each iteration the likelihood is evaluated on

764     a different subset of data.  The '-lowmem' flag has ChromHMM reduce main memory usage by not storing

765     in main memory all the input data and instead re-loading from disk when needed.

766

767     **Choice of number of states**

768          We trained full-stack models with 20, 40, 60, 80, 100 and 120 states, using the data and procedure

769     outlined above. We then quantitatively compared the chromatin state annotations from these models in

770     terms of their power to predict locations of various other genomic annotations not used in the model

771     training: Exon, Gene Body, TSS, TSS2kb, CpG Islands, TES, laminB1lads elements (listed in section

772    *External Annotation Sources* section). Specifically, we evaluated the predictive power using the AUROCs

773    that are calculated as described in a subsection below. Across different genomic contexts, as the number of

774    full-stack states increased, the AUROC increased, but the marginal increase was smaller as the number of

775    states increased (**Supp. Fig. 3**). To balance the additional information available in models with increased

776    number of states, while keeping the number of states manageable for interpretation and downstream

777    analysis, we choose to focus on a model with 100 states. We note that this choice is greater than previously

778    used for cell-type-specific chromatin state models [3,16,21], reflecting the additional information available

779    for genome annotation based on the large number of datasets spanning many cell types that we are using.

780

781    **Lifting chromatin state annotations to hg38**

782        The chromatin state annotation resulted from stacked modeling was in hg19. In order to obtain the

783    annotations for hg38, we first wrote the chromatin state map hg19 in .bed format such that each line

784    corresponds to a genomic region of 200bp. We then used liftOver tools downloaded from UCSC utilities

785    to generate the chromatin state annotation in hg38. In total, there are 1,186,379 200-bp segments that were

786    not mapped from hg19 to hg38.

787

788    **Summary sets of experiments**

789        To construct a summary visualization of the emission parameters with a reduced set of features that

790    approximate the annotation from the full model, we applied a greedy search over the 1032 input datasets as

791    described in **Supplementary Methods.** We applied this procedure to reduce the 1032 input datasets to 80

792    summary datasets.

793

794    **Identifying states with differential association of marks for individual tissue groups**

795        For each state, we tested for combinations of the 8 most profiled marks, and 19 tissue groups

796    previously defined [16], whether the emission probabilities of features associated with one chromatin mark

797 and in one tissue group was significantly greater than those of features associated with the same mark and

798 not in the tissue group. The eight marks that we tested were H3K9me3, H3K4me1, H3K4me3, H3K27me3,

799 H3K36me3, H3K27ac, H3K9ac, and DNase. H3K27ac, H3K9ac and DNase were profiled in 98, 62 and 53

800 reference epigenomes, respectively, and the remaining five marks in 127 reference epigenomes. For tests

801 involving H3K27ac, H3K9ac, and DNase, we excluded tissue groups for which there were no experiments.

802 In total, there were 14,200 tests among 100 states, 8 chromatin marks and 19 tissue groups. For each

803 combination of state, chromatin mark and tissue group being tested, we applied a one-sided Mann-Whitney

804 test to test whether the emission probabilities of the state for the features associated with the tested mark in

805 the tested tissue group are greater than those in other tissue groups. The Bonferroni-corrected p-value

806 threshold based on a significance level of 0.05 to declare a test significant was 3.5e-6.

807

808 **Computing coefficients of variation across different tissue groups**

809 For each state, we looked into the emission probabilities of experiments associated with six

810 chromatin marks strongly associated with promoter and enhancer activities (DNase, H3K27ac, H3K4me1,

811 H3K4me2, H3K4me3, H3K9ac). We grouped these experiments based on their associated chromatin mark

812 and tissue groups, and calculated the average emission probabilities of experiments in each chromatin mark-

813 tissue group combination. For each state and chromatin mark combination, we then calculated the

814 coefficient of variation across different tissue groups, in terms of average emission probabilities from the

815 previous step. For each group of states, we averaged the resulting coefficients of variation across states of

816 the same group. The results show the average coefficients of variation of emission probabilities across

817 different tissue groups for each state group- chromatin mark combination.

818

819 **Computing fold enrichments for other annotations**

820 All overlap enrichments for external annotations were computed using the ChromHMM

821 OverlapEnrichment command. We used the '-b 1' flag, which specifies a binning resolution of the

822 annotations. This '-b 1' flag is necessary when computing enrichments based on the hg38 liftOver

36

823    annotations, which no longer respects the 200bp segment coordinate intervals from hg19. Including this

824    flag gives the same results when applied to annotations from hg19 with 200bp segments, though with extra

825    computational costs. We also included the '-lowmem' flag to specify the lower memory usage option. The

826    ChromHMM command OverlapEnrichment computes fold enrichment between chromatin states and

827    provided external annotations relative to a uniform genome-wide background distribution. More

828    specifically, the fold enrichments are calculated as:

829
$$FE_{x,s} = \frac{\frac{\#SX}{\#X}}{\frac{\#S}{\#G}} = \frac{\frac{\#SX}{\#S}}{\frac{\#X}{\#G}} = \frac{\#SX \cdot \#G}{\#S \cdot \#X}$$

830    where

831    $FE_{x,s}$: fold enrichment of state $s$ in genomic context $x$

832    $\#S$: number of genomic positions belonging to the state $S$

833    $\#X$: number of genomic positions where genomic context $X$ is present

834    $\#SX$: number of genomic bins that overlap both state $S$ and genomic context $X$

835    $\#G$: number of genomic positions in the entire genome

836

837    **Enrichment with cell-type-specific ChromHMM annotations**

838        We computed the enrichments of the full-stack states for cell-type-specific ChromHMM chromatin

839    state annotations. For the cell-type-specific chromatin state annotations we used 25-state ChromHMM

840    annotations of 127 reference epigenomes from the Roadmap Epigenomics project. This model was trained

841    using the concatenated modeling approach using imputed data of 12 chromatin marks [16,22]. For each of

842    the 100 full-stack states, we calculated the enrichment for the 25 states separately in each of the 127

843    reference epigenomes, resulting in 127 tables of 25 enrichment values for each of the 100 states. We

844    summarized this information by reporting, for each of the 100 full-stack states, and 127 reference

845    epigenomes, the cell-type-specific state among the 25 states that is maximally enriched, resulting in a 100-

846    by-127 table. We also summarized the information by reporting for each of the 100 full-stack states and 25

847     cell-type-specific states, the maximum and median fold enrichments across the 127 reference epigenomes

848     (**Supplementary Data**).

849

850     **Receiver operator characteristic curve analysis for predicting external annotations**

851     To evaluate the information available in the chromatin state annotations from a chromatin state

852     model that can help predict locations of an external genomic annotation, we computed the Receiver

853     Operator Characteristic (ROC). To do this, we first divided the genome into 200bp bins, and randomly

854     partitioned 50% of the bins for training and the remaining 50% for testing. For a target external genome

855     annotation, we computed the enrichment of such annotation with each chromatin state on the training data.

856     We then ranked states in decreasing order of enrichments for the target annotation. We used this ranking of

857     states to iteratively add genomic bases assigned to the added state to our predictions of bases overlapping

858     the target annotation in the testing dataset. Based on the overlap of the predictions and the target annotation

859     at each iteration, we plotted ROC curves and summarized the information by computing area under the

860     ROC curves (AUROC).

861

862     **Cell-type-specific ChromHMM annotations for comparing predictive information**

863     We compared the full-stack model to two sets of cell-type-specific annotations in terms of their

864     ability to predict external annotations. One set of cell-type-specific annotations was the 18-state

865     ChromHMM from Roadmap Epigenomic Project [16], which was trained using observed data for six

866     chromatin marks: H3K4me1, H3K4me3, H3K9me3, H3K27ac, H3K27me3 and H3K36me3, using the

867     concatenated approach.

868     The second set of cell-type-specific ChromHMM annotations were annotations we generated here

869     to have a more stringent comparison. We partitioned the 1032 datasets we used to learn the full-stack model

870     into 127 subsets based on their associated reference epigenome. For each of these 127 subsets, we applied

871     ChromHMM to learn a cell-type-specific model with 100 states. We learned these models with the same

872     procedure as described above for the full-stack model, with the exception of using the '-init random' flag

873     to randomly initialize models' parameters. This flag was necessary since for some reference epigenomes,

874     the number of specified states (100) was greater than the number of combinations of input datasets, which

875     is the maximum number of states supported by ChromHMM default initialization. We specified the number

876     of states as 100 in these cell-type-specific models to control for the number of states in comparing with the

877     full-stack model.  However, we note that due to the large number of states relative to the input tracks, some

878     of these models ended up having fewer than 100 states being assigned to positions in the genome.

879

880     **Computing fine-mapped variant enrichment**

881          To compute enrichment of full-stack states for phenotypically associated fine-mapped variants, we

882     downloaded data on fine-mapped variants for 3052 traits from CAUSALdb [52]. Specifically we obtained

883     posterior probabilities of variants being causal based on two fine-mapping methods, FINEMAP [54] and

884     CAVIAR [53], which do not use epigenomic annotations as part of the fine mapping procedure. For each

885     method and trait combination, we separately partitioned the provided set of potential causal variants into

886     distinct loci. To form the distinct loci, we merged neighboring variants into the same loci until there was at

887     least 1MB-gap between the two closest variants from different loci. Separately for each fine-mapping

888     method, trait, and locus combination, we selected the single variant with the highest posterior probability

889     of being causal. For each fine-mapping method, we took the union of variants across 3052 traits, and then

890     calculated the fold enrichments for the union of these lead variants with stacked ChromHMM states relative

891     to the enrichment with a background set of common variants from dbSNP build 151 (hg19). To do this, we

892     separately computed the enrichments of both of these sets relative to a genome-wide background, and then

893     divided the enrichment of the foreground set (lead fine-mapped variants) by the enrichment of the

894     background set (common variants). The dbSNP variants were obtained from the UCSC genome browser.

895

896     **Computing structural variant enrichments**

897          To compute enrichment of the full-stack states for structural variant enrichments, we obtained data

898     of structural variants from [33]. We used the B38 call set, which was in hg38 and used for the analysis

899     presented in [33]. We filtered out structural variants that did not pass the quality control criteria of [33]. We

900     then separately considered structural variants annotated as either a deletion or a duplication, for which there

901     were, 112,328 and 28,962 sites respectively.

902     Since the structural variants were defined in hg38, we computed their enrichment for ChromHMM

903     state annotations in full-stack and cell-type-specific models that were lifted over from hg19 to hg38,

904     following the procedure outlined above. Next, we followed the enrichment analysis procedure outlined

905     above to compare full-stack vs. cell-type-specific chromatin state segmentations' power in recovering

906     structural variants.

907     To compare the power of full-stack state annotations vs. cell-type-specific state annotation

908     frequency, we utilized 15-state genome-wide chromatin state data for 127 cell types (reference epigenomes)

909     from Roadmap Epigenomics Consortium. We followed the analysis outlined in [33], for each of the 15

910     ChromHMM states, we annotated genomic positions based on the number of cell types in which the state

911     is present (ranging from 0 to 127), resulting in 15 state-specific models' annotations. We then applied the

912     procedure above to compare the predictive power of different models' annotations against the full-stack

913     annotation. For the state-specific models, the enrichment values are calculated for structural variants and

914     number of cell types that a ChromHMM state is assigned to.

915

916     **Computing enrichments with cancer-associated variants**

917     We obtained data of somatic mutations associated with different types of cancer from COSMIC

918     non-coding variants dataset v.88 in hg38 [55]. We selected from this dataset variants that were from whole-

919     genome sequencing. We filtered out variants that overlap with any of the following: the hg38 black-listed

920     regions from the ENCODE Data Analysis Center (DAC) [58], hg38 dbSNP (v151) set of common variants

921     from the UCSC genome browser database, or regions annotated as coding sequence ('CDS') based on

922     GENCODE v.30 hg38 [59] gene annotations. We decided to restrict this analysis to the four cancer types

923     with the most number of variants present in the dataset in hg38: liver (1,351,417), pancreas (500,930),

924     haematopoietic and lymphoid tissue (354,501), and breast (323,751), we then lifted over these sets of

925    variants from hg38 to hg19, resulting in 1,351,159, 500,798, 354,351, and 323,685, variants respectively.

926    To obtain a background set of genomic locations for the enrichment analysis, we filtered from the genome

927    the same set of hg38 annotations of black-listed regions, common variants, and coding sequences. We then

928    lifted over these remaining positions from hg38 to hg19 to obtain the background. We calculated the

929    enrichment of chromatin states with cancer-associated variants by first calculating the enrichment values

930    of chromatin states with filtered variants associated with each of the four cancer types, and the enrichment

931    values with background set of genomic bases, all relative to the whole genome. We then divided the cancer-

932    associated variant enrichment values by the background bases enrichments.

933

**934    External annotations sources**

935    The sources for external annotations for enrichments analyses, not given above, were as follows:

936       • CpG island annotations were those included in the ChromHMM (v1.18) and originally obtained

937           from the UCSC genome browser.

938       • Annotations of exon, gene bodies, transcription start (TSS), and transcription end sites (TES),

939           2kb windows surrounding TSSs (TSS2kb) were RefSeq annotations included in ChromHMM

940           (v1.18) and originally based on annotations obtained from the UCSC genome browser.

941       • Lamina associated domains were for human embryonic lung fibroblasts that were included in

942           ChromHMM (1.18), which were lifted over to hg19 from hg18 positions originally provided

943           by [60].

944       • Annotations of assembly gaps were obtained from the UCSC genome browser and correspond

945           to the Gap track.

946       • Annotations of zinc finger (ZNF) genes correspond to coordinates of genes whose name

947           contained 'ZNF' from GENCODE's hg19 gene annotation, v30 [59].

948       • Annotations of coding sequences correspond to coordinates of genes whose feature type is

949           'CDS' from GENCODE's hg19 gene annotation, v30 [59].

950     •    Annotations of pseudogenes correspond to coordinates of genes those whose gene type or

951          transcript type contained 'pseudogene' from GENCODE's hg19 gene annotation, v30 [59].

952     •    Annotations of repeat elements were obtained from UCSC genome browser RepeatMasker

953          hg19 tracks.

954     •    Cell-type-specific ChromHMM chromatin state annotations were obtained from the Roadmap

955          Epigenomics Consortium through http://compbio.mit.edu/roadmap [16]. These include data of

956          the 15-state and 18-state models based on observed data and the 25-state chromatin model

957          based on imputed data for 127, 98 and 127 reference epigenomes, respectively.

958     •    CTCF- cell-type-specific chromatin states were based on the ChromHMM chromatin state

959          annotations for six human cell types (GM12878, H1ESC, Helas3, Hepg2, Huvec, K562) for a

960          25-state model from the ENCODE integrative analysis [22,30]. We extracted coordinates of

961          region annotated to the 'Ctcf' and 'CtcfO', both associated with CTCF signal and limited

962          histone mark signal.

963     •    Blacklisted regions were those provided by the ENCODE Data Analysis Center (DAC) for

964          hg19 and hg38 [58].

965     •    ConsHMM conservation state annotations for human (hg19) were those from [32].

966     •    Annotations of human genetic variants and their allele frequency were from GNOMAD v2.1.1

967          [47]. The dataset includes 229 million SNVs and 33 million indels from 15,708 genomes of

968          unrelated individuals, which are aligned against the GRCg37/hg19 reference.

969     •    GWAS catalog variants were obtained from the NHGRI-EBI Catalog, accessed on December

970          5th, 2016 [50].

971

972

973     **Analysis of gene expression across states**

974    To analyze the relationship between gene expression and the full-stack states, we downloaded gene

975    expression data from the Roadmap Epigenomics Consortium [16]. Specifically, we downloaded a matrix

976    of gene expression values, in RPKM (Reads Per Kilobase Million), for protein coding genes for 56 reference

977    epigenomes that were among the 127 used as part of the full-stack model. In total, we obtained expression

978    values for 19,795 Ensembl protein coding genes.

979    The gene expression data was obtained from

980    (https://egg2.wustl.edu/roadmap/data/byDataType/rna/expression/57epigenomes.exon.RPKM.pc.gz).  We

981    also obtained the corresponding genomic coordinates for these genes from

982    (https://egg2.wustl.edu/roadmap/data/byDataType/rna/expression/Ensembl_v65.Gencode_v10.ENSG.gen

983    e_info). For this analysis, we filtered out genes that are not classified as protein-coding. We transformed

984    the gene expression values by adding a pseudo-count of 1 to the raw counts in RPKM, and taking the log

985    of the resulting values.

986    For each full-stack-state and 56 reference epigenomes, we calculated the average gene expression

987    of all genes overlapping with the state, taking into account the genes' length. For each gene $g$ we denote

988    its length $L_g$ and expression $E_g$. We let $s_i$ denote the state assigned at the 200-bp bin $i$ and $G_i$ denote the

989    set of genes overlapping the 200bp bin $i$. Let $B_s$ denote the set of 200bp bins that are assigned to state $s$.

990    The average normalized expression with state $s$ then becomes:

991    $$avg \exp bp \ normalized_s = \frac{\sum_{i \in B_s} \sum_{g \in G_i} E_g / L_g}{\sum_{i \in B_s} \sum_{g \in G_i} 1 / L_g}$$

992    We also calculated for each state the average and coefficient of variation of these averages across

993    reference epigenomes. We used the BEDTools *bedtools intersect* command to obtain the chromatin state

994    assignments for 200bp segments that totally or partially overlap with any gene. To obtain average gene

995    expressions of a state in a cell group as presented in **Fig. 3C**, we averaged the reported bp-normalized

996    average gene expressions of the corresponding state across cell types within the group.

997    We also analyzed average gene expression values for each state as a function of the position of the

998    state annotations relative to TSS, following a procedure similar to what was used previously [3]. We first

43

999     identified a gene's outer transcription start site (TSS) based on the reported coordinates of the gene and

1000    strand in the gene annotation file noted above. For each 200bp bin that is within 25kb upstream or

1001    downstream of an annotated TSS, including those that directly overlap with an annotated TSS, we

1002    determined the assigned full-stack state at this bin, and the position of the bin relative to those TSSs. Bins

1003    directly overlapping an annotated TSS were at position 0. If the gene was on the positive strand, the

1004    segments' genomic coordinates lower than the TSSs' correspond to upstream regions at negative points

1005    (minimum value: -250000), while genomic coordinates higher than the TSSs' correspond to downstream

1006    regions at positive points (maximum value: 25000). If the gene is on the negative strand, the upstream and

1007    downstream positions are reversed. For each state and each 200-bp bin position relative to TSS, we

1008    determined the subset of genes where there is a 200bp bin annotated to that state at that position relative to

1009    their TSSs, and calculated their average expression. This produces a 100-by-251 table for one reference

1010    epigenome, corresponding to the number of full-stack states and 200-bp segments intersecting the 50kb

1011    windows surrounding genes' TSSs and one segment directly overlapping the TSSs. We then smoothened

1012    the averaged expression data spatially by applying the sliding window average algorithm with a window

1013    size of 21, i.e. each segment's smoothened gene expression is the average of data in that segment and 21

1014    surrounding genomic segments. Data of average gene expression in the first and last 10 segments within

1015    the 50kb window are not included in the window of smoothened data. We averaged results of 56 tables

1016    corresponding to 56 reference epigenomes as the final output from this procedure.

1017

1018    **Computing enrichment for bases prioritized by variant prioritization scores**

1019            To compute state enrichments for bases prioritized by different variant prioritization scores, we

1020    followed the approach of [32]. We obtained coordinates of bases containing prioritized variants based on

1021    14 different methods as processed and described in [32]. The scores were Eigen and Eigen-PC version 1.1,

1022    funSeq2 version 2.1.6, and CADD v1.4, REMM, FIRE, fitCons, CDTS, LINSIGHT, FATHMM, GERP++,

1023    phastCons, phyloP and DANN [34–46]. For 12 of the 14 scores, we separately considered prioritized

1024    variants genome-wide and in non-coding regions only. Two of the variant prioritization scores, LINSIGHT

44

1025 and FunSeq2 [36,38], were defined only in the non-coding regions, so these scores were only used in the

1026 non-coding region analysis. As described in [32], the regions included in the non-coding analysis were

1027 defined as the bases where both LINSIGHT and FunSeq2 provided scores, which was 90.4% of the genome.

1028 For both the non-coding and whole genome analysis we computed the enrichment for bases ranked in the

1029 top 1%, 5% or 10% using the variant prioritization scores. We note that because of ties in some scores, the

1030 score-threshold above which we classified the bases as prioritized was chosen to be as close as possible to

1031 the target percentage (1%, 5% or 10%). We also note that if there were any bases with missing values for

1032 any particular score, then that base was assigned with the minimum values of such scores.

1033 Enrichment values for the whole genome were computed as described above with the

1034 OverlapEnrichment command from ChromHMM. For computing enrichments restricted to non-coding

1035 regions, we first calculated enrichment of the non-coding prioritized variants relative to the whole genome

1036 and the enrichment of non-coding regions as defined above relative to the whole genome. We then divided

1037 these two enrichment values to obtain the enrichment of prioritized non-coding variants within non-coding

1038 regions.

1039

1040 **Data availability**

1041 Full-stack chromatin state annotation of the genome is available at

1042 https://github.com/ernstlab/full_stack_ChromHMM_annotations. An updated version of ChromHMM is

1043 available at https://ernstlab.biolchem.ucla.edu/ChromHMM/

1044

1045

1046 **Acknowledgements**

1051    Jonsson Comprehensive Cancer Center and Eli and Edythe Broad Center of Regenerative Medicine and

1052    Stem Cell Research Ablon Scholars Program.

1053

1054    **Ethics Declarations**

1055    The authors announce no conflicts of interests.

1056     **References**

1057     1. Barski A, Cuddapah S, Cui K, Roh T-Y, Schones DE, Wang Z, et al. High-resolution profiling

1058     of histone methylations in the human genome. Cell. Elsevier; 2007;129:823–837.

1059     2. Boyle AP, Davis S, Shulha HP, Meltzer P, Margulies EH, Weng Z, et al. High-resolution

1060     mapping and characterization of open chromatin across the genome. Cell. Elsevier;

1061     2008;132:311–322.

1062     3. Ernst J, Kheradpour P, Mikkelsen TS, Shoresh N, Ward LD, Epstein CB, et al. Mapping and

1063     analysis of chromatin state dynamics in nine human cell types. Nature. Nature Publishing Group;

1064     2011;473:43–49.

1065     4. Thurman RE, Rynes E, Humbert R, Vierstra J, Maurano MT, Haugen E, et al. The accessible

1066     chromatin landscape of the human genome. Nature. Nature Publishing Group; 2012;489:75–82.

1067     5. Xie W, Schultz MD, Lister R, Hou Z, Rajagopal N, Ray P, et al. Epigenomic analysis of

1068     multilineage differentiation of human embryonic stem cells. Cell. Elsevier; 2013;153:1134–

1069     1148.

1070     6. Taberlay PC, Statham AL, Kelly TK, Clark SJ, Jones PA. Reconfiguration of nucleosome-

1071     depleted regions at distal regulatory elements accompanies DNA methylation of enhancers and

1072     insulators in cancer. Genome Res. Cold Spring Harbor Lab; 2014;24:1421–1432.

1073     7. Claussnitzer M, Dankel SN, Kim K-H, Quon G, Meuleman W, Haugen C, et al. FTO obesity

1074     variant circuitry and adipocyte browning in humans. N Engl J Med. Mass Medical Soc;

1075     2015;373:895–907.

1076    8. Gjoneska E, Pfenning AR, Mathys H, Quon G, Kundaje A, Tsai L-H, et al. Conserved

1077    epigenomic signals in mice and humans reveal immune basis of Alzheimer's disease. Nature.

1078    Nature Publishing Group; 2015;518:365–369.

1079    9. Kheradpour P, Ernst J, Melnikov A, Rogov P, Wang L, Zhang X, et al. Systematic dissection

1080    of regulatory motifs in 2000 predicted human enhancers using a massively parallel reporter

1081    assay. Genome Res. Cold Spring Harbor Lab; 2013;23:800–811.

1082    10. Varshney A, Scott LJ, Welch RP, Erdos MR, Chines PS, Narisu N, et al. Genetic regulatory

1083    signatures underlying islet gene expression and type 2 diabetes. Proc Natl Acad Sci. National

1084    Acad Sciences; 2017;114:2301–2306.

1085    11. Lay FD, Triche TJ, Tsai YC, Su S-F, Martin SE, Daneshmand S, et al. Reprogramming of

1086    the human intestinal epigenome by surgical tissue transposition. Genome Res. Cold Spring

1087    Harbor Lab; 2014;24:545–553.

1088    12. Lee J, Krivega I, Dale RK, Dean A. The LDB1 complex co-opts CTCF for erythroid lineage-

1089    specific long-range enhancer interactions. Cell Rep. Elsevier; 2017;19:2490–2502.

1090    13. Consortium EP. Identification and analysis of functional elements in 1% of the human

1091    genome by the ENCODE pilot project. nature. Nature Publishing Group; 2007;447:799.

1092    14. Consortium EP. An integrated encyclopedia of DNA elements in the human genome. Nature.

1093    Nature Publishing Group; 2012;489:57–74.

1094    15. Fernández AF, Bayón GF, Urdinguio RG, Toraño EG, García MG, Carella A, et al.

1095    H3K4me1 marks DNA regions hypomethylated during aging in human stem and differentiated

1096    cells. Genome Res. Cold Spring Harbor Lab; 2015;25:27–40.

1097    16. Kundaje A, Meuleman W, Ernst J, Bilenky M, Yen A, Heravi-Moussavi A, et al. Integrative

1098    analysis of 111 reference human epigenomes. Nature. Nature Publishing Group; 2015;518:317–

1099    330.

1100    17. Mikkelsen TS, Ku M, Jaffe DB, Issac B, Lieberman E, Giannoukos G, et al. Genome-wide

1101    maps of chromatin state in pluripotent and lineage-committed cells. Nature. Nature Publishing

1102    Group; 2007;448:553–560.

1103    18. Wang Q, Yu G, Ming X, Xia W, Xu X, Zhang Y, et al. Imprecise DNMT1 activity coupled

1104    with neighbor-guided correction enables robust yet flexible epigenetic inheritance. Nat Genet.

1105    Nature Publishing Group; 2020;52:828–839.

1106    19. Zhu J, Adli M, Zou JY, Verstappen G, Coyne M, Zhang X, et al. Genome-wide chromatin

1107    state transitions associated with developmental and environmental cues. Cell. Elsevier;

1108    2013;152:642–654.

1109    20. Stunnenberg HG, Abrignani S, Adams D, de Almeida M, Altucci L, Amin V, et al. The

1110    International Human Epigenome Consortium: a blueprint for scientific collaboration and

1111    discovery. Cell. Elsevier; 2016;167:1145–1149.

1112    21. Ernst J, Kellis M. Discovery and characterization of chromatin states for systematic

1113    annotation of the human genome. Nat Biotechnol. Nature Publishing Group; 2010;28:817–825.

1114    22. Ernst J, Kellis M. ChromHMM: automating chromatin-state discovery and characterization.

1115    Nat Methods. Nature Publishing Group; 2012;9:215–216.

1116    23. Hoffman MM, Buske OJ, Wang J, Weng Z, Bilmes JA, Noble WS. Unsupervised pattern

1117    discovery in human chromatin structure through genomic segmentation. Nat Methods. Nature

1118    Publishing Group; 2012;9:473.

1119    24. Libbrecht MW, Rodriguez OL, Weng Z, Bilmes JA, Hoffman MM, Noble WS. A unified

1120    encyclopedia of human functional DNA elements through fully automated annotation of 164

1121    human cell types. Genome Biol. Springer; 2019;20:180.

1122    25. Ernst J, Kellis M. Chromatin-state discovery and genome annotation with ChromHMM. Nat

1123    Protoc. Nature Publishing Group; 2017;12:2478.

1124    26. Biesinger J, Wang Y, Xie X. Discovering and mapping chromatin states using a tree hidden

1125    Markov model. BMC Bioinformatics. Springer; 2013. p. S4.

1126    27. Zhang Y, An L, Yue F, Hardison RC. Jointly characterizing epigenetic dynamics across

1127    multiple human cell types. Nucleic Acids Res. Oxford University Press; 2016;44:6721–6731.

1128    28. Chronis C, Fiziev P, Papp B, Butz S, Bonora G, Sabri S, et al. Cooperative binding of

1129    transcription factors orchestrates reprogramming. Cell. Elsevier; 2017;168:442–459.

1130    29. Mortazavi A, Pepke S, Jansen C, Marinov GK, Ernst J, Kellis M, et al. Integrating and

1131    mining the chromatin landscape of cell-type specificity using self-organizing maps. Genome

1132    Res. Cold Spring Harbor Lab; 2013;23:2136–2148.

1133    30. Hoffman MM, Ernst J, Wilder SP, Kundaje A, Harris RS, Libbrecht M, et al. Integrative

1134    annotation of chromatin elements from ENCODE data. Nucleic Acids Res. Oxford University

1135    Press; 2013;41:827–841.

1136    31. Heintzman ND, Stuart RK, Hon G, Fu Y, Ching CW, Hawkins RD, et al. Distinct and

1137    predictive chromatin signatures of transcriptional promoters and enhancers in the human

1138    genome. Nat Genet. Nature Publishing Group; 2007;39:311–318.

1139    32. Arneson A, Ernst J. Systematic discovery of conservation states for single-nucleotide

1140    annotation of the human genome. Commun Biol. Nature Publishing Group; 2019;2:1–14.

1141    33. Abel HJ, Larson DE, Regier AA, Chiang C, Das I, Kanchi KL, et al. Mapping and

1142    characterization of structural variation in 17,795 human genomes. Nature. Nature Publishing

1143    Group; 2020;583:83–89.

1144    34. Cooper GM, Goode DL, Ng SB, Sidow A, Bamshad MJ, Shendure J, et al. Single-nucleotide

1145    evolutionary constraint scores highlight disease-causing mutations. Nat Methods. Nature

1146    Publishing Group; 2010;7:250–251.

1147    35. Di Iulio J, Bartha I, Wong EH, Yu H-C, Lavrenko V, Yang D, et al. The human noncoding

1148    genome defined by genetic diversity. Nat Genet. Nature Publishing Group; 2018;50:333–337.

1149    36. Fu Y, Liu Z, Lou S, Bedford J, Mu XJ, Yip KY, et al. FunSeq2: a framework for prioritizing

1150    noncoding regulatory variants in cancer. Genome Biol. BioMed Central; 2014;15:1–15.

1151 37. Gulko B, Hubisz MJ, Gronau I, Siepel A. A method for calculating probabilities of fitness

1152 consequences for point mutations across the human genome. Nat Genet. Nature Publishing

1153 Group; 2015;47:276–283.

1154 38. Huang Y-F, Gulko B, Siepel A. Fast, scalable prediction of deleterious noncoding variants

1155 from functional and population genomic data. Nat Genet. Nature Publishing Group;

1156 2017;49:618–624.

1157 39. Ioannidis NM, Davis JR, DeGorter MK, Larson NB, McDonnell SK, French AJ, et al. FIRE:

1158 functional inference of genetic variants that regulate gene expression. Bioinformatics. Oxford

1159 University Press; 2017;33:3895–3901.

1160 40. Ionita-Laza I, McCallum K, Xu B, Buxbaum JD. A spectral approach integrating functional

1161 genomic annotations for coding and noncoding variants. Nat Genet. Nature Publishing Group;

1162 2016;48:214.

1163 41. Pollard KS, Hubisz MJ, Rosenbloom KR, Siepel A. Detection of nonneutral substitution rates

1164 on mammalian phylogenies. Genome Res. Cold Spring Harbor Lab; 2010;20:110–121.

1165 42. Quang D, Chen Y, Xie X. DANN: a deep learning approach for annotating the pathogenicity

1166 of genetic variants. Bioinformatics. Oxford University Press; 2015;31:761–763.

1167 43. Rentzsch P, Witten D, Cooper GM, Shendure J, Kircher M. CADD: predicting the

1168 deleteriousness of variants throughout the human genome. Nucleic Acids Res. Oxford University

1169 Press; 2019;47:D886–D894.

1170    44. Rogers MF, Shihab HA, Mort M, Cooper DN, Gaunt TR, Campbell C. FATHMM-XF:

1171    accurate prediction of pathogenic point mutations via extended features. Bioinformatics. Oxford

1172    University Press; 2018;34:511–513.

1173    45. Siepel A, Bejerano G, Pedersen JS, Hinrichs AS, Hou M, Rosenbloom K, et al.

1174    Evolutionarily conserved elements in vertebrate, insect, worm, and yeast genomes. Genome Res.

1175    Cold Spring Harbor Lab; 2005;15:1034–1050.

1176    46. Smedley D, Schubach M, Jacobsen JO, Köhler S, Zemojtel T, Spielmann M, et al. A whole-

1177    genome analysis framework for effective identification of pathogenic regulatory variants in

1178    Mendelian disease. Am J Hum Genet. Elsevier; 2016;99:595–606.

1179    47. Karczewski KJ, Francioli LC, Tiao G, Cummings BB, Alföldi J, Wang Q, et al. The

1180    mutational constraint spectrum quantified from variation in 141,456 humans. Nature. Nature

1181    Publishing Group; 2020;581:434–43.

1182    48. Francioli LC, Polak PP, Koren A, Menelaou A, Chun S, Renkens I, et al. Genome-wide

1183    patterns and properties of de novo mutations in humans. Nat Genet. Nature Publishing Group;

1184    2015;47:822–826.

1185    49. Lindblad-Toh K, Garber M, Zuk O, Lin MF, Parker BJ, Washietl S, et al. A high-resolution

1186    map of human evolutionary constraint using 29 mammals. Nature. Nature Publishing Group;

1187    2011;478:476–482.

1188    50. Welter D, MacArthur J, Morales J, Burdett T, Hall P, Junkins H, et al. The NHGRI GWAS

1189    Catalog, a curated resource of SNP-trait associations. Nucleic Acids Res. Oxford University

1190    Press; 2014;42:D1001–D1006.

1191    51. Hindorff LA, Sethupathy P, Junkins HA, Ramos EM, Mehta JP, Collins FS, et al. Potential

1192    etiologic and functional implications of genome-wide association loci for human diseases and

1193    traits. Proc Natl Acad Sci. National Acad Sciences; 2009;106:9362–9367.

1194    52. Wang J, Huang D, Zhou Y, Yao H, Liu H, Zhai S, et al. CAUSALdb: a database for

1195    disease/trait causal variants identified using summary statistics of genome-wide association

1196    studies. Nucleic Acids Res. Oxford University Press; 2020;48:D807–D816.

1197    53. Chen W, Larrabee BR, Ovsyannikova IG, Kennedy RB, Haralambieva IH, Poland GA, et al.

1198    Fine mapping causal variants with an approximate Bayesian method using marginal test

1199    statistics. Genetics. Genetics Soc America; 2015;200:719–736.

1200    54. Benner C, Spencer CC, Havulinna AS, Salomaa V, Ripatti S, Pirinen M. FINEMAP:

1201    efficient variable selection using summary data from genome-wide association studies.

1202    Bioinformatics. Oxford University Press; 2016;32:1493–1501.

1203    55. Tate JG, Bamford S, Jubb HC, Sondka Z, Beare DM, Bindal N, et al. COSMIC: the

1204    catalogue of somatic mutations in cancer. Nucleic Acids Res. Oxford University Press;

1205    2019;47:D941–7.

1206    56. Schuster-Böckler B, Lehner B. Chromatin organization is a major influence on regional

1207    mutation rates in human cancer cells. nature. Nature Publishing Group; 2012;488:504–507.

1208    57. Parker SC, Gartner J, Cardenas-Navia I, Wei X, Abaan HO, Ajay SS, et al. Mutational

1209    signatures of de-differentiation in functional non-coding regions of melanoma genomes. PLoS

1210    Genet. Public Library of Science; 2012;8:e1002871.

1211    58. Amemiya HM, Kundaje A, Boyle AP. The ENCODE blacklist: identification of problematic

1212    regions of the genome. Sci Rep. Nature Publishing Group; 2019;9:1–5.


1213    59. Harrow J, Frankish A, Gonzalez JM, Tapanari E, Diekhans M, Kokocinski F, et al.

1214    GENCODE: the reference human genome annotation for The ENCODE Project. Genome Res.

1215    Cold Spring Harbor Lab; 2012;22:1760–74.


1216    60. Guelen L, Pagie L, Brasset E, Meuleman W, Faza MB, Talhout W, et al. Domain

1217    organization of human chromosomes revealed by mapping of nuclear lamina interactions.

1218    Nature. Nature Publishing Group; 2008;453:948–951.

1219

1220

1221

1222

1223

1224

1225