

1 **Supplementary Information**

2 **Methods**

3 **Quantitative metric of model selection**

4 To determine a model of chromatin states that most closely represented our thyroid
5 derived data, we selected a model with the most discrete and inter-sample consistent
6 output state emissions. In other words, the model that is the most well defined,
7 maximizing the homogeneity of epigenetic features in chromatin states across samples.
8 Concretely, choosing a model based on this selection metric will make it so that the set
9 of epigenetic features associated with a region partitioned as state 2 in one sample will
10 tend to be similar (or homogeneous) to the set of epigenetic features associated with a
11 region partitioned as state 2 in another sample. We provide an R package (hmpickr
12 available at <https://github.com/csiu/hmpickr>) to help users select such a model
13 (doi:10.5281/zenodo.398681). Overall, we choose the model that has the lowest
14 homogeneity cost, which we compute as follows:

15

16 Let H represent the total number of histone marks and h represent a particular histone
17 mark. Here $h = \{1, 2, \dots, H\}$. We represent a probability close to 0 or 1, representing
18 respectively absent and present histone marks across regions of the same state, by
19 taking the minimum of the emission probability of a histone mark for a state (E_{hk}) and 1
20 minus that probability. To increase the penalty on states that are not as well defined,
21 state costs are squared. To account for the difference in the number of states across
22 models, we normalize by the number of states (K) in each model. Overall, we represent

23 the homogeneity cost of the state (d_k) and the homogeneity cost of a model (D) as
24 follows:

$$d_k = \sum_{h=1}^H \min\{1 - E_{hk}, E_{hk}\}$$
$$D = \frac{\sum_{k=1}^K d_k^2}{K}$$

25

26 **Determination of Chromatin states**

27 We used ChromHMM v1.12 (Ernst & Kellis 2012), an implementation of a hidden
28 Markov model, to learn combinatorial chromatin states jointly across 8 thyroid
29 epigenomes (a normal and diseased thyroid sample from each of the 4 thyroid sample
30 donors). ChromHMM was trained using 6 histone marks (H3K4me1, H3K4me3,
31 H3K27ac, H3K36me3, H3K9me3, and H3K27me3). For each ChIP-seq data set, read
32 counts were computed in non-overlapping 200bp bins across the entire genome. In total
33 there were 15,181,508 bins. Each bin was discretized using ChromHMM's BinarizeBam
34 into two levels: 1 indicating enrichment, and 0 indicating no enrichment. The binarization
35 was performed by comparing ChIP-seq read counts to ChIP-seq input control data for
36 local adjustments to the binarization threshold. We have also used ChIP-seq input
37 (obtained before immunoprecipitation) control data as an additional input feature directly
38 in the model. Reads mapping to chromosome Y were discarded to ensure reads that
39 were mismapped were not carried forward in the computation. Command "LearnModel"
40 with options "-p 11" was specified to use 11 processors in parallel to train a model using
41 a standard Baum-Welch training algorithm. We trained a total of 26 models with the
42 number of states ranging from 11 to 23 states. The trained model was then used to

43 compute the posterior probability of each state for each genomic bin in each sample.
44 The regions were labelled using the state with the maximum posterior probability. To
45 assign biologically meaningful labels to the states, we used ChromHMM package to
46 compute the overlap and neighborhood enrichments of each state relative to
47 coordinates of known functional annotation obtained from the Epigenome Roadmap
48 Project (Roadmap Epigenomics Consortium *et al.* 2015). The chromatin state models
49 and browser tracks can be downloaded from <http://www.bcgsc.ca/data/thyroid>.

50

51 **Estimating transcript abundance, gene expression, and gene variance**

52 We used Salmon v0.7.2 (Patro *et al.* 2017) to estimate transcript abundance from RNA-
53 seq reads. As input, Salmon takes a reference transcriptome and a set of raw sequence
54 reads. Each read is 75nt in length. The transcriptome was downloaded from the UCSC
55 Table Browser with options as follows: group “Genes and Gene Predictions”, track
56 “GENCODE Genes V19”, table “Basic (wgEncodeGencodeBasicV19)”, and output
57 format “sequence”. The function “salmon index” was used to index the reference
58 transcriptome, while “salmon quant” was used to estimate transcript abundance
59 measured in transcripts per million (TPM). To integrate the transcript-level abundance
60 estimates into gene-level abundance estimates, the tximport function of the tximport R
61 package v1.2.0 (Soneson *et al.* 2015) was used to sum up the Salmon estimated
62 transcript abundances within genes. This was also repeated for read counts. The
63 regularized logarithm transformation (rlog) function of the DESeq2 R package v1.14.0
64 (Love *et al.* 2014) was then used to transform tximport generated read count data to

65 render them homoskedastic (i.e. such that the variance of the errors over the samples
66 are similar). Gene variance was calculated on the rlog transformed read counts.

67

68 **Selecting genes that have low expression in non-thyroid tissue types**

69 Gene expression of various tissue types were obtained from the Genotype-Tissue
70 Expression (GTEx) project. Data was downloaded for “Query: Genes matching: “,
71 specifically expressed in any Organism part above the expression level cutoff: 0 in
72 experiment E-MTAB-2919” at <https://www.ebi.ac.uk/gxa/experiments/E-MTAB-2919> on
73 November 21, 2016. Expression is measured in Fragments Per Kilobase of transcript
74 per Million mapped reads (FPKM). We consider a gene as lowly expressed if the FPKM
75 is less than or equal to 10. Expression values were then binarized to “low” and “high”
76 expression. Genes for 52 non-thyroid samples were then clustered and visualized on a
77 heat map. The cluster of genes that had low expression across all non-thyroid samples
78 were then considered the set of genes that had low expression in 52 non-thyroid tissue
79 types.

80 **Motif analysis of Active Enhancers nearby genes**

81 Sequence analysis of the genomic DNA in regions marked as active enhancer (state
82 10) by all four samples nearby (within 2 Mbp) the highly expressed 42 unique protein
83 coding genes indicated 9 transcription factor motifs (p-value < 0.05). These transcription
84 factor motifs (and consensus sequence) are Fosl2 (NATGASTCABNN), Jun-AP1
85 (GATGASTCATCN), Fra1 (NNATGASTCATH), MafK (GCTGASTCAGCA), BATF
86 (DATGASTCAT), Atf3 (DATGASTCATHN), RORgt (AAYTAGGTCA), ZSCAN22
87 (SMCAGTCWGAKGGAGGAGGC), and Bach2 (TGCTGAGTCA). After Benjamini
88 multiple test correction, the first 3 motifs (i.e. Fosl2, Jun-AP1, and Fra1) were found to
89 be significantly enriched.

90

91 A similar sequence analysis for the 10 highly expressed genes consistent across the
92 four specimens indicated 0 transcription factor motifs (p-value < 0.05).

93

94 Furthermore, a similar sequence analysis for the 18 genes that we consider
95 epigenetically active and consistently expressed in the thyroid that are likely highly
96 relevant to thyroid function indicated 2 transcription factor motifs (p-value < 0.05). These
97 transcription factor motifs (and consensus sequence) are LXRE
98 (RGGTTACTANAGGTCA) and ZNF675 (ARGAGGMCAAATGW). After Benjamini
99 multiple test correction, no motifs were found to be significantly enriched.

100

101 The list of all motifs and their significance are found in the Supplementary Excel file.

102 **Distinct patterns of repressive marks**

103 Detailed methodology for bisulfite-seq and data processing is available in the
104 Supplemental Experimental Procedures of (Pellacani et al. 2016), at
105 <http://www.epigenomes.ca/protocols-and-standards>, or upon request. We used FindER
106 v1.0.0b (available at <http://www.epigenomes.ca/tools-and-software/finder/index.html>)
107 with default options to find enriched ChIP-seq regions and bedtools v2.24.0 (Quinlan *et*
108 *al.* 2010) to integrate DNA methylation profiles with repressive chromatin marks.
109 Integration was done by (1) finding a list of H3K4me3 or H3K27me3 repressed regions
110 unique to each individual, (2) mapping the fractional methylation call per 200 bp
111 genomic bin onto the list of repressed regions, and (3) comparing the DNA methylation
112 profiles across individuals for each of the repressed regions.

113

114

115

116

117

118

119

120

121

122

123 Quinlan AR & Hall IM 2010 BEDTools: a flexible suite of utilities for comparing genomic
124 features. *Bioinformatics* **26** 841–842