

Dynamics of gene expression and chromatin marking during cell state transition

Beatrice Borsari¹, Amaya Abad¹, Cecilia C. Klein^{1,2,3}, Ramil Nurtdinov¹, Alexandre Esteban^{1,4}, Emilio Palumbo¹, Marina Ruiz-Romero¹, María Sanz^{1,5}, Bruna R. Correa¹, Rory Johnson^{1,6}, Sílvia Pérez-Lluch^{1,*} and Roderic Guigó^{1,7,*}

¹Centre for Genomic Regulation (CRG), The Barcelona Institute of Science and Technology, Barcelona 08003, Catalonia, Spain

²Departament de Genètica, Microbiologia i Estadística, Facultat de Biologia and Institut de Biomedicina (IBUB), Universitat de Barcelona, Barcelona 08028, Catalonia, Spain

³Present address: Clarivate, Barcelona 08025, Catalonia, Spain

⁴Present address: "la Caixa" Foundation, Department of Research and Innovation, Barcelona 08028, Catalonia, Spain

⁵Present address: Universidad Camilo José Cela (UCJC), Madrid 28692, Spain

⁶Present address: Department of Medical Oncology, Inselspital, University Hospital and University of Bern, Bern 3010, Switzerland. School of Biology & Environmental Science, University College Dublin (UCD), Dublin 4, Ireland.

⁷Universitat Pompeu Fabra (UPF), Barcelona 08003, Catalonia, Spain

*Correspondence: silvia.perez@crg.cat (S.P-L.), roderic.guigo@crg.cat (R.G.)

1 **Abstract**

2 A large body of data strongly supports a crucial role for histone modifications in the regulation of gene
3 expression. An increasing number of cases, however, are being reported in which changes in gene expres-
4 sion occur without changes in histone modifications. To provide a framework where to properly investigate
5 these apparently contradictory observations, we have generated an unprecedented time series deep epige-
6 nomics data during a transdifferentiation process that occurs with massive transcriptional changes. During
7 this process, we find a strong coupling between histone modifications and gene expression only at the time
8 of initial gene activation, when deposition of marks, mostly following gene activation, tends to occur in a
9 precise order. Other than at that time, changes in gene expression are mostly uncoupled from changes
10 in histone modifications. Over genes, these occur in a very limited number of combinations, defining the
11 major chromatin states in the genome, which are largely stable. The overall association between gene
12 expression and chromatin is, thus, much weaker than reported so far in steady-state conditions and, for
13 some marks, actually runs in the opposite direction.

14 Introduction

15 Chromatin is the complex of DNA, histone and non-histone proteins that constitutes the chromosomes
16 found in the nucleus of eukaryotic cells. Post-translational modifications (PTMs) of histone proteins, to-
17 gether with other epigenetic features, can alter the overall chromatin structure and are thought to play a
18 critical role in the regulation of all DNA-based processes. In particular, interest has grown in understanding
19 the relationship between chromatin and transcriptional regulation.

20 Histone marks have been assumed to play an important role in the regulation of gene expression, asso-
21 ciated with either active or silent gene expression. For instance, high levels of H3K27ac and H3K4me1 are
22 considered a feature of active transcriptional enhancers¹, whereas active promoters are typically marked by
23 H3K4me3^{2,3}. Conversely, constitutive and facultative heterochromatin is normally associated with higher
24 levels of H3K9me3 and H3K27me3, respectively^{4,5}. According to the histone code hypothesis⁶, distinct
25 combinations of histone modifications over regulatory regions — associated with specific arrangement of
26 transcription factors — confer to each gene a unique temporal and spatial transcriptional program. In
27 strong support of this hypothesis, methods to predict gene expression from combinations of different hi-
28 stone marks have been developed with great accuracy, even when the predictions are obtained in a cell
29 type other than the one in which the model is inferred^{7,8}. This presumed association between gene ex-
30 pression and histone modifications underlines, for instance, the great amount of efforts invested to target
31 chromatin modifications, with the goal of altering gene expression, to treat certain diseases such as cancer
32 and degenerative diseases^{9–11}.

33 The majority of these predictions are conducted in steady-state conditions, and therefore do not track
34 the association between gene expression and histone marks over time. Studies along time, however, are
35 essential to decipher the mechanisms behind transcriptional control and maintenance, since an appro-
36 priate balance of stability and dynamics in epigenetic features seems to be required for accurate gene
37 expression¹². Interestingly, a number of studies in different species and biological models have highlighted
38 a degree of correlation between gene expression and chromatin marks over time substantially lower (or
39 even absent) than what previously described in steady-state conditions. For instance, during fruit fly devel-
40 opment, around 34% of the expressed genes lack H3K4me3 at their promoters¹³, while transcription can
41 occur in the absence of most active marks^{14,15}. It has also been reported that, upon stimulation, changes
42 in gene expression are not always accompanied by changes in histone modifications¹⁶, and that chromatin
43 marks do not represent linear measures of transcriptional activity^{17,18}. Overall, it has been suggested
44 that the contribution of chromatin to gene expression may partially depend on the promoter architecture of
45 genes¹⁹.

46 Time-series studies have also striven to elucidate the temporal ordering in which transcription factor
47 (TF) binding, deposition of histone marks and RNA Polymerase recruitment occur at both enhancer and
48 promoter regions. For instance, it has been reported that enhancers required for hematopoietic differ-
49 entiation are already primed with H3K4me1 in multipotent progenitors²⁰. However, *de novo* enhancers'
50 transcription seems to precede local deposition of H3K4me1 and H3K4me2 marks²¹. Furthermore, de-

51 position of H3K4me1 is dispensable for either enhancer or promoter transcription, and does not affect the
52 maintenance of transcriptional programs^{22,23}.

53 Nevertheless, most time-series studies so far have monitored a few histone modifications in a limited
54 number of time-points. To address these limitations, here we have generated gene expression profiles and
55 maps of nine histone modifications at twelve time-points along a controlled cellular differentiation process:
56 the induced transdifferentiation of human BLaER1 cells into macrophages²⁴. BLaER1 is a human B-cell
57 precursor leukemia cell line, stably transfected with a construct containing cEBP α fused with the estrogen
58 hormone receptor binding domain²⁴. These cells are able to transdifferentiate into functional macrophages
59 at a high efficiency rate upon induction with beta-estradiol, which induces the internalization of the tran-
60 scription factor into the nucleus, promoting massive transcriptomic changes. We believe that the data that
61 we have generated constitutes an unprecedented resource in the field to understand epigenetic regulation
62 of gene expression.

63 Analysis of these data reveals that the large steady-state associations between gene expression and
64 chromatin marking previously reported are partially artifactual, and mainly arise from the constrained nature
65 of the transcriptome and the epigenome. When measured over time, these correlations are globally weak
66 and, remarkably, in the case of H3K9me3, run in the opposite direction that previously thought. We found
67 that, in contrast to the histone code hypothesis, only a limited number of combinations of histone modifi-
68 cations are actually marking the genes, defining the major genic chromatin states in the human genome.
69 Genes tend to remain in the same state throughout the entire transdifferentiation process, even those that
70 change expression substantially. We have also observed substantial chromatin changes that are not nec-
71 essarily accompanied by changes in gene expression, suggesting that epigenetic modifications contribute
72 to cell state in a manner that cannot be fully recapitulated by gene expression. We did find, however, a strong
73 association between chromatin marking and expression at the time of initial gene activation. We have been
74 able to determine the precise order of histone modifications at that time, and found that only H3K4me1 and
75 H3K4me2 appear to be deposited prior to gene activation. Further changes in gene expression, compara-
76 ble or even stronger than those at gene activation, seem to be mostly uncoupled from changes in histone
77 modifications.

78 **A rich resource for time-series analysis of chromatin and gene expression dynamics**

79 To investigate the temporal interplay between transcriptional activity and chromatin marking during the
80 transdifferentiation of BLaER1 cells into macrophages²⁴, we monitored this process at 12 time-points, from
81 0 to 168 hours post-induction (p.i.) (Figure 1a). Reciprocal regulation of B-cell and macrophage antigens
82 CD19 and Mac-1, respectively, was assessed by flow cytometry throughout the process (Supplementary
83 Figure 1a).

84 For each time-point we characterized, in two biological replicates, the whole cell RNA-seq gene ex-
85 pression profiles and the ChIP-seq maps of nine histone post-translational modifications. Besides the
86 six marks (H3K4me1, H3K4me3, H3K27ac, H3K27me3, H3K36me3 and H3K9me3) endorsed by the ref-

87 erence epigenome criteria (International Human Epigenome Consortium, <http://ihec-epigenomes.org/research/reference-epigenome-standards/>), we have profiled H3K4me2, H3K9ac and H4K20me1
88 (Figure 1b). This has allowed us to characterize the interchange between different degrees of lysine four
89 methylation over time, but also to compare acetylation patterns on distinct lysine residues, and to explore
90 the alternation of broad marks over actively transcribed gene bodies. In addition, we have generated, for
91 each time-point, ChIP-seq profiles of the transcription factor cEBP α , RNA-seq data from the cytosol and
92 the nucleus, as well as riboprofiling and proteomics maps (Correa *et al.*, in preparation).

93
94 To avoid any bias due to differences in the transdifferentiation process between experiments, a crucial
95 component of our experimental design is that the RNA and the chromatin to perform immunoprecipitations
96 with all histone marks were obtained from the same pool of cells in each biological replicate (see Methods).
97 To efficiently and reproducibly analyze the wealth of data generated in a controlled environment, we de-
98 veloped *ChIP-nf* (<https://github.com/guigolab/chip-nf>), a pipeline implemented in NextFlow²⁵
99 (see Methods).

100 **Gene expression recapitulates transdifferentiation more accurately than chromatin**

101 To characterize gene expression and histone modifications' profiles during the pre-B cell transdifferentiation
102 process, we selected the 12,248 genes — out of 19,831 protein-coding genes annotated in Gencode²⁶ ver-
103 sion 24 — that were either expressed in at least one time-point (≥ 5 TPM, 10,696 genes), or silent all along
104 the process (0 TPM in all time-points, 1,552 genes) (Supplementary Figure 1b). Within expressed genes,
105 we identified 8,030 genes characterized by significant changes in their expression profiles over time (differ-
106 entially expressed, DE; Supplementary Figure 1b; see Methods). Half of these genes are down-regulated
107 during the process, 25% are up-regulated, and for the remaining 25% we observed transient increases
108 (peaking) and decreases (bending) in expression. 2,666 expressed genes do not display changes in ex-
109 pression over time (stably expressed).

110 For every gene in these sets, we also computed the level of each histone modification at a specific
111 time-point, either over the gene body in the case of H3K36me3 and H4K20me1, or at promoter regions
112 (± 2 Kb with respect to the transcription start site) for the remaining marks (Supplementary Figure 1c, see
113 Methods). Roughly all expressed genes are marked by the canonical active histone modifications, whereas
114 the proportion of silent genes showing peaks of these marks is low, except for H3K4me1 and H3K4me2
115 (Supplementary Table 1). Unexpectedly, marks typically associated with silent transcription (H3K9me3 and
116 H3K27me3) are not abundant in either expressed or silent genes.

117 To visually summarize the gene expression and individual histone modification profiles during transd-
118 ifferentiation, we performed Principal Component Analysis (PCA), in which we plotted the 12 time-points
119 based on these profiles (Figure 1c). Even though the PCA was performed jointly on gene expression and
120 all chromatin marks — which show different patterns of variation —, the first two principal components
121 (PC1 and PC2) still capture about one fifth of the total variance of the data. Whereas gene expression is
122 able to recapitulate the process in the space of the first two principal components, the chromatin marks are

123 less resolute, with H3K27ac, H3K9ac and H4K20me1 showing the clearest trends. The trajectory of gene
124 expression in the PCA space suggests that the process occurs in two different transcriptional phases, with
125 PC1 explaining the main differences between pre-B cells and macrophages, and PC2 representing early
126 transcriptional changes within the first 24 hours of transdifferentiation. Instead, for several chromatin marks
127 we observed parabolic trajectories, with PC2 mainly separating the intermediate stages of transdifferentia-
128 tion from the differentiated cell types. Genes contributing to PC1 are mostly up- or down-regulated (Sup-
129 plementary Figure 1d), and display significant enrichment in Gene Ontology terms associated with immune
130 response and cell motility (Supplementary Table 2). Instead, PC2-contributing genes perform functions re-
131 lated to nucleic acids metabolism and protein modification (Supplementary Table 2), and comprise a large
132 proportion of genes either displaying no changes in gene expression, or presenting transient increases or
133 decreases (Supplementary Figure 1d). Taken all together, these results suggest that, while there are major
134 changes in gene expression and chromatin leading from one differentiated cell type to another (PC1), there
135 are also changes that may be involved in a transient de-differentiation from pre-B cells into an intermediate
136 state, and in the re-differentiation into macrophages (PC2), with expression contributing differently from
137 chromatin marks.

138 **The association between chromatin marking and gene expression is overestimated by cor-** 139 **relations computed in steady-state conditions**

140 We computed, at each time-point, the steady-state correlation between levels of expression and histone
141 modifications across the set of 12,248 genes (Figure 1d). As previously observed, we found a strong
142 positive correlation for most active marks (median Pearson r value across time-points between 0.51 and
143 0.72), and a (weak) negative correlation for the repressive marks H3K9me3 and H3K27me3 (-0.07 and
144 -0.17, respectively). However, when computing, for individual genes, the correlation between expression
145 and chromatin profiles through time (time-course correlations), the values are substantially lower for active
146 marks (median Pearson r ranging between 0.10 and 0.45), and higher for repressive marks (0.13 and -
147 0.03 for H3K9me3 and H3K27me3, respectively; Figure 1d). Remarkably, for H3K9me3 the time-course
148 correlation with expression is positive, in contrast to the repressive role generally assumed for this mark
149 (see, for instance,²⁷).

150 It appears, therefore, that correlations measured in steady-state conditions artificially inflate the true
151 degree of association between gene expression and chromatin modifications, and even mis-represent the
152 direction of this association. This can be dramatically seen by randomizing the real temporal associa-
153 tion between gene expression and chromatin marks. Within each gene's time-series profile, we permuted
154 histone modification levels among time-points, while keeping the actual gene expression values (see Meth-
155 ods; for an example with H3K4me3, compare upper and lower panels in Supplementary Figure 2a). As
156 expected, the average time-course correlation is zero for all marks (Supplementary Figure 2b). However,
157 the steady-state correlations are unexpectedly large for canonically active marks upon randomization, de-
158 spite the fact that any meaningful association between gene expression and chromatin marks has been

159 eliminated (Supplementary Figures 2a lower panel and 2b). This is likely due to a considerable fraction of
160 genes displaying stable expression and chromatin profiles over time, which are either relatively highly ex-
161 pressed and marked (housekeeping genes)²⁸, or silent and not marked. Indeed, after removing the genes
162 with silent or stable expression profiles over time, the steady-state correlations (Supplementary Figure 2c)
163 are lower compared to those computed on the entire set of genes (Figure 1d), and become more similar to
164 the time-course correlations.

165 **Genes are characterized by a limited number of major chromatin states, which are more** 166 **stable than expression**

167 Next, we investigated the dynamics of chromatin marking during transdifferentiation. Towards that end, we
168 summarized the chromatin state of each gene at each time-point, by building a multivariate Hidden Markov
169 Model (HMM) on the signal of the nine histone marks along the twelve transdifferentiation points. More
170 specifically, we produced a segmentation of the transdifferentiation time by assigning a given chromatin
171 state to each gene at each time-point. This is in contrast to previous uses of HMMs in the field, where
172 the segmentation is produced along the genome sequence by assigning a given chromatin state to every
173 genome interval^{29–33}. We explored configurations with up to twenty different states, and found that five
174 states are a good compromise between optimizing the likelihood of the model and the number of states
175 capturing the epigenetic status of genes (Supplementary Figure 3a and Figure 2a, see Methods). These
176 five states correspond to the major combinations of histone modifications in which genes can be found
177 (major chromatin states): a) absence of marking, with the exception of moderate H3K9me3 signal, b) low
178 marking (mono and di-methylation of H3K4), c) bivalent marking (mostly marking by H3K4me1, H3K4me2
179 and/or H3K27me3), d) canonical active marking (all canonical active marks) and e) strong canonical active
180 marking in the presence of H3K9me3 signal. These states (from a to e) correspond to increasing marking
181 by canonically active histone modifications, with the exception of the bivalent marking state (c), which is
182 also characterized by high H3K27me3 signal. These results suggest that only a limited number of combi-
183 nations of marks can co-occur in a given gene at a given time-point. They also indicate that marking by
184 H3K4me1 and H3K4me2 appears to be a precondition for marking by any other active histone modification,
185 since for none of the configurations that we have explored, we have found states in which there is mark-
186 ing by an active histone modification without H3K4me1 and H3K4me2. The most frequent states among
187 expressed genes are active and strong active marking (d and e, respectively), while the most frequent
188 state among silent genes is absence of marking (a) (Supplementary Figure 3c). This state is defined by
189 moderate marking by H3K9me3, consistent with the assumed repressive role of this mark. Strong marking
190 by H3K9me3, however, defines also unexpectedly the strong marking state (e), characteristic of expressed
191 genes (Supplementary Figure 3b). This, together with the overall positive temporal correlation of this mark
192 with gene expression (Figure 1d), suggests a so far unappreciated dual role for this mark. Indeed, we have
193 found this mark both over genes silent along transdifferentiation (Supplementary Figure 3c), as well as over
194 up-regulated (Fig. 2d, middle panel) and stably expressed genes (Supplementary Figure 3d).

195 Hierarchical clustering of genes based on the sequence of the five states along the twelve time-points
196 revealed a limited number of temporal chromatin state profiles (Figures 2b-c). Most of the genes remain in
197 the same chromatin state during transdifferentiation (constant state profiles), irrespective of whether they
198 are stably (79%) or differentially expressed (70%) along the process (Figure 2d, left panels). Thus, during
199 transdifferentiation, most changes in gene expression are not accompanied by chromatin changes. Of the
200 remaining genes, the vast majority (90%) go over just one-state transition during transdifferentiation. When
201 considering DE genes, these transitions are generally associated with the expected transcriptional changes
202 (Figure 2c). Transitions from weaker to stronger active chromatin marking are accompanied by increases
203 in gene expression (Figure 2c, upper side; Figure 2d, middle panels), while transitions from stronger to
204 weaker active chromatin states are accompanied by decreases in gene expression (Figure 2c, lower side;
205 Figure 2d, right panels). However, while transitions from active to strong active marking states (and vice
206 versa) are more numerous, the corresponding fold changes in gene expression are lower, compared to
207 transitions from low marking to active marking states (and vice versa). We observed activating transitions
208 from the absent state mainly to the low marking state, further supporting the fact that marking by H3K4me1
209 and H3K4me2 is a prerequisite for the deposition of any other active histone modifications. On the other
210 hand, we did not observe transitions from the strong active marking state to absence of marking, suggesting
211 that the erasing of chromatin marks is not as an efficient process as its deposition.

212 Analysis of individual histone marks confirmed the HMM results. We determined whether the marks'
213 signals are stable or variable over time, analogously to what was done for gene expression profiles. The
214 majority of genes present, indeed, stable chromatin profiles during transdifferentiation, even when focusing
215 only on the differentially expressed ones (Supplementary Table 3, left side; Figure 3a). Lysine acetylation
216 (H3K27ac and H3K9ac) is the most dynamic signal (Supplementary Table 3, left side). Still, around 35%
217 of DE genes show no changes in histone acetylation, despite being marked. Unexpectedly, only 8.5% of
218 DE genes show changes in H3K27me3 throughout the process, although roughly half of them are down-
219 regulated. Conversely, for a smaller number of silent and stably expressed genes we observed significant
220 variations in their chromatin profiles over time (Supplementary Table 4, Figures 3b-c), comparable or even
221 larger than for DE genes (Supplementary Figure 4a), although no changes could be detected in their
222 expression profiles.

223 We observed, in general, that differentially marked genes display clearer transdifferentiation trajectories
224 compared to genes that are stably marked (Supplementary Figure 4b), further supporting that the contri-
225 bution of gene expression and chromatin marks to cell state is not fully overlapping. Consistent with the
226 positive association between H3K9me3 and gene expression, the trajectory for this mark resembles more
227 the trajectories of some active marks such as H3K4me1 and H3K4me2, than that of H3K27me3. Actu-
228 ally, we have also found more genes in which H3K9me3 is positively than negatively correlated with gene
229 expression (see Supplementary Table 3, right side).

Chromatin marking is associated with expression specifically at the time of gene activation

The limited number of chromatin HMM states indicates a coordinated behaviour of histone modifications. To investigate this behaviour at the resolution of individual marks and how it relates to gene expression, we first determined the type of association between each mark and expression along transdifferentiation, for each of the 8,030 genes that are differentially expressed (labels: unmarked, stably marked, positively correlated, uncorrelated and negatively correlated; see Figure 4a, Supplementary Table 3 and Methods). Then, we clustered the combinations of marks and types of association, and found that, in general, in a given gene, most marks show indeed the same type of association with expression (Figure 4b). When clustering the genes based on these combinations, we found essentially three major groups (Figure 4c, Supplementary Figure 5a). The first and largest cluster includes 4,995 DE genes (62%), presenting either stable or uncorrelated profiles for the majority of active marks, and absence of marking for H3K27me3 and H3K9me3 (Figures 5a-b, upper panels). The second cluster includes 2,993 DE genes (37%), showing the canonical positive correlation between expression and most active modifications. A large proportion of these genes lack repressive marks, but a few of them (9%) exhibit the expected negative correlation with H3K27me3 (Figures 5a-b, middle panels). Finally, the third and smallest cluster includes 102 genes (1%) characterized by an overall absence of both active and repressive marking, with the exception of H3K4me1 and H3K4me2 (Figures 5a-b, lower panels).

Especially in the case of up-regulated genes, these clusters mostly reflect the level of gene activation when transdifferentiation starts (Figure 5c, Supplementary Figures 5b-c). Genes in cluster 1 are already activated at the beginning of transdifferentiation, genes in cluster 2 are in early stages of activation or are activated early during transdifferentiation, while genes in cluster 3 are activated late during the process. The functions of the genes in these clusters are consistent with their level of activation at the beginning of transdifferentiation (Supplementary Figures 5d-e). In particular, genes in cluster 3 are associated with macrophage-specific functions, and we have found them lowly expressed and lowly marked in other cell types but CD14⁺ monocytes (Supplementary Figures 5f-g). Down-regulation of gene expression, on the other hand, appears to be largely uncoupled from chromatin changes, since most genes decreasing expression belong to cluster 1 (Supplementary Figure 5h).

Gene expression changes anticipate changes in most active marks for up-regulated genes

The results above are suggestive that the association between gene expression and histone modifications occurs preferentially in a limited window of time during the initial stage of gene activation. Thus, to investigate the relationship between expression and chromatin marking precisely at this stage, we focused on the set of 257 up-regulated genes that are not expressed at 0 hours p.i., and that are, therefore, specifically activated during transdifferentiation. The vast majority of these genes (230, 89%) belong to cluster 2, that is, they are indeed characterized by positive correlation between gene expression and active chromatin marks. They are mostly associated with low and bivalent marking HMM states and, in 25% of the cases,

265 transition into stronger marking states towards the end of transdifferentiation (Supplementary Figure 6a,
266 upper panel).

267 To investigate the temporal relationship between gene activation and chromatin marking, for each up-
268 regulated gene and histone mark we rescaled the expression and chromatin time-series profiles to the
269 same range (0-100%), and identified the first time-point at which the expression level and the chromatin
270 signal reach at least 25%, 50%, 75% and 100% (Supplementary Figure 6b). In this way, we determined
271 whether active chromatin marking anticipates, co-occurs with, or follows gene expression. In contrast to
272 the prevalent view, we did not find that most active marks anticipate activation of gene expression. At the
273 first stage of up-regulation (25%), only marking by H3K4me1, H3K4me2 and H3K27ac anticipates more
274 often than follows activation of gene expression (Figures 6a-b), whereas for the other marks most changes
275 follow expression up-regulation. These differences are progressively lost towards the end of the process
276 (Figure 6a, Supplementary Figure 6c).

277 To further decipher the precise order in which active chromatin signals are established over time, we
278 computed, for a given mark, the fraction of genes whose changes either anticipate (Figure 6c, upper panel)
279 or co-occur with (Supplementary Figure 6d, upper panel) changes in each of the other six marks. When
280 considering 25% of up-regulation, we observed that, in general, no marks anticipate H3K4me1, indicating
281 that it is the first mark to increase, followed by H3K4me2 and H3K27ac (Figure 6c, upper panel). This is
282 consistent with the HMM analysis, which suggested that marking by H3K4me1 and H3K4me2 is a pre-
283 requisite for marking by other histone modifications (Figure 2a). Changes in H3K4me1, H3K4me2 and
284 H3K27ac most frequently precede increases in H3K9ac and H3K4me3. In all the comparisons, H3K36me3
285 and H4K20me1 follow the other marks (Figure 6c, upper panel). As observed for gene expression, this
286 precise order of marks' deposition is progressively lost along transdifferentiation (Figure 6c upper panel,
287 Supplementary Figure 6d upper panel). Overall, this suggests that the deposition of active chromatin mod-
288 ifications follows a precise order at the time of initial gene activation (H3K4me1 > H3K4me2 > H3K27ac >
289 expression > H3K9ac > H3K4me3 > H3K36me3 > H4K20me1; Figure 6d, left panel).

290 We performed a similar analysis with the set of 629 up-regulated genes that are already substantially
291 expressed at 0 hours p.i. (> 25 TPM). These genes belong mostly to cluster 1 (389, 62%), that is, their
292 expression profiles are uncoupled from changes in chromatin marking, and they actually remain in active
293 chromatin states during transdifferentiation (Supplementary Figure 6a lower panel). For these genes we
294 did not find preservation in the pattern of chromatin deposition with respect to expression (Supplementary
295 Figure 6e), nor in the deposition of the marks (Figure 6c lower panel; Figure 6d right panel; Supplementary
296 Figure 6d lower panel).

297 **A model to explain the coupling between transcription and chromatin marking over time**

298 Altogether, our results show that the canonical association between histone modifications and gene ex-
299 pression mainly occurs in a limited window of time preceding and following initial gene activation. We
300 specifically propose a model (Figure 7a) in which the activation of gene expression is anticipated by de-

301 position of H3K4me1, H3K4me2 and, less frequently, of H3K27ac at promoter regions. The deposition of
302 other marks typically enriched either at promoters (H3K9ac, H3K4me3) or over the gene body (H3K36me3,
303 H4K20me1) is concomitant to or, more often, follows (and may be induced by) gene activation. After this
304 initial stage of gene activation, further changes in gene expression, comparable or even stronger, appear
305 to be mostly uncoupled from changes in histone modifications (Figure 7b, compare left and right panels).

306 This model explains our observations well. The patterns of association between chromatin marking
307 and gene expression (as defined in Figure 4a) for genes in different degrees of activation when transdif-
308 ferentiation starts (0h p.i.) reflect how this association changes as gene activation proceeds (Figure 7c).
309 Up-regulated genes that are silent when transdifferentiation starts (mostly in cluster 3) lack almost all “ac-
310 tivating” histone modifications, possibly with the exception of H3K4me1 and H3K4me2 (i.). Up-regulated
311 genes in cluster 2 that are lowly or not activated at 0h show mostly correlated patterns of expression and
312 chromatin marking. In these genes, most marks, with the exception of H3K4me1, H3K4me2 and H3K27ac,
313 follow rather than anticipate expression (ii., see also Figure 7b, left panel). As we consider genes with
314 increasing degrees of activation at 0h (and thus, in increasingly advanced states of activation), the fraction
315 of genes with correlated patterns of expression and chromatin marking decreases, while the fraction of
316 genes with stable or uncorrelated chromatin profiles (iii. and iv.) proportionally increases. The temporal
317 order of activation of marks observed in early activation stages is also gradually lost. Finally, for genes in
318 cluster 1 (v.), which are already highly active when transdifferentiation starts, changes in gene expression,
319 even if substantial, are mostly uncoupled from chromatin marking, showing uncorrelated or stable profiles
320 (see also Figure 7b, right panel).

321 Discussion

322 Epigenetics was initially defined as “the branch of biology that studies the causal interactions between
323 genes and their products which bring the phenotype into being”³⁴. In a more contemporary definition, “an
324 epigenetic trait is a stably heritable phenotype resulting from changes in a chromosome without alterations
325 in the DNA sequence”³⁵. The epigenetic mechanisms leading to the development of an individual or
326 to the differentiation of a cell lineage from the unique genotype of the organism have been largely studied
327 during decades. Although initial references to the mechanisms by which epigenetics promotes cell memory
328 and leads cell fate did not relate to its ability to regulate gene expression, a causative role for epigenetic
329 modifications in controlling transcription has been later pointed out (see 36,37 for reviews about different
330 aspects related to epigenetics and its role in regulating gene expression), and it has even been shown that
331 some epigenetic features, such as histone modifications, are accurate predictors of gene expression^{7,8,38}
332 and the other way around³⁹.

333 However, the causal/consequential relationship between chromatin modifications and gene expression
334 represents a long-standing discussion^{40,41}, and a number of reports have challenged the causal role that
335 has been broadly attributed to chromatin modifications^{14,22,42,43}. Still, and despite the efforts dedicated to

336 this problem and the vast literature produced, the actual relationship between histone modifications and
337 the regulation of gene expression remains unsolved.

338 This is partially due to the few available studies in which gene expression and histone modifications have
339 been both consistently monitored through time in a given dynamic system. Differentiation models are suit-
340 able to study the relationship between gene expression and chromatin marking, as they provide a dynamic
341 system that allows to decipher the order of the events. In this work, we have used the transdifferentiation of
342 BLaER1 cells (pre-B cells) into macrophages, a model that has proven to be highly efficient²⁴, and we have
343 generated high-quality data on the transcriptome and the epigenome in twelve time-points along the seven
344 days the transdifferentiation process lasts. Our analysis of these data has uncovered some fundamental
345 features of chromatin organization in human genes and of the relationship between gene expression and
346 histone modifications.

347 Our analyses have also contributed to a better understanding of the molecular events underlying trans-
348 differentiation of pre-B cells into macrophages. Despite the fact that, to our knowledge, there is no retro-
349 differentiation during the process^{24,44}, the joint PCA of gene expression and chromatin marks suggests
350 that BLaER1 cells undergo an intermediate state (Figure 1c). This intermediate state is characterized by
351 chromatin changes not accompanied by changes in gene expression (Supplementary Figure 7), and vice
352 versa by changes in gene expression not associated with chromatin changes (Supplementary Figure 7a).
353 Although it is often assumed that the transcriptome is the main determinant of cell state, these results
354 suggest that epigenetic modifications contribute to cell state in a manner that cannot be fully recapitulated by
355 gene expression. Thus, neither the epigenome nor the transcriptome can be fully predictive of one another.

356 Consistently, we found that the association between gene expression and chromatin modifications is
357 overall weaker than reflected by the correlations reported so far, which have been mostly computed in a
358 particular steady-state cellular condition (Figure 1d). These artifactually strong correlations result from the
359 largely constrained nature of the human epigenome and transcriptome. In particular, a large fraction of
360 genes in the human genome (likely more than 50%²⁸) are either invariably silent and not marked, or ex-
361 pressed and marked across most cellular states. Genes with stable epigenomes and transcriptomes drive
362 the correlations to large values when computed in a particular cell condition, and explain why models re-
363 lating gene expression to histone modifications inferred in a particular cell type have high predictive power
364 in other cell types^{7,8,38,39}, even though there is no true causality involved in the relationship between chro-
365 matin and expression. The steady-state correlations represent an example of the Simpson's paradox⁴⁵, by
366 which the data can show different or even opposite behavior if subgroups within the dataset are considered.

367 HMMs have been widely used to summarize patterns of combinations of multiple histone modifications
368 into a limited number of chromatin states. However, in most cases so far, they have been used to segment
369 the genome sequence²⁹⁻³³. Here, instead, we used them, we believe for the first time, to segment time
370 along a dynamic differentiation process. The HMM segmentation reveals that, even though the number
371 of possible histone combinations is very large (if nine histones are considered, $2^9 = 512$ combinations
372 are possible), most genes are actually found in one among only about five major states (Figure 2a). This

373 challenges to some extent the notion of a histone code⁶. Further supporting the limited number of genic
374 chromatin states, we found that marks act in a coordinated manner, meaning that genes showing a stable
375 profile for one histone modification tend also to present stable profiles of the other marks, and that genes
376 showing absence of one active mark tend to be void of all positive modifications (Figures 4b-c, Supplemen-
377 tary Figure 5a). Most genes remain in the same chromatin state during transdifferentiation, irrespective of
378 whether they are or not differentially expressed, explaining the low correlation between gene expression
379 and chromatin marks throughout time. Analysis of individual histone modifications confirmed these obser-
380 vations, and further identified a number of silent or stably expressed genes along transdifferentiation that
381 show changes in chromatin marking (Figures 3b-c).

382 While we have not extensively focused on marks typically associated with gene silencing, our analy-
383 ses have nevertheless uncovered some unexpected findings regarding these marks. First, we observed
384 that, although roughly 4,000 genes are down-regulated during the process, only 10% of them present
385 H3K27me3 marking in at least one time-point, indicating that the majority of genes that are silenced along
386 transdifferentiation do not depend on Polycomb repression. Most remarkably, however, we have found that
387 H3K9me3 is actually more often associated with gene activation than with gene silencing, in contrast to
388 what has been previously reported²⁷. While H3K9me3 at the transcription start site has been previously
389 related to active expression in malignant cells⁴⁶ and, more recently, to actively transcribed genes in early
390 preimplantation embryos⁴⁷, our results show that H3K9me3 is likely to have a general dual association,
391 both with up- and down-regulation of gene expression. Additional analyses are required to understand the
392 conditions under which H3K9me3 plays either role, but our HMM suggest that H3K9me3 alone is associ-
393 ated with repression, while when acting in conjunction with other marks is positively associated with gene
394 expression.

395 While there is a general lack of coupling between gene expression and chromatin marking, there is
396 a temporal relationship between gene expression and the different histone modifications at the time of
397 gene activation. We propose a model (Figure 7a) in which activation of gene expression is anticipated by
398 deposition of H3K4me1, H3K4me2, while deposition of other marks is concomitant or, more often, follows
399 gene activation, being the gene body marks the last ones to be incorporated. The order of chromatin
400 marking in our model is in agreement with the observed deposition of histone modifications upon induction
401 of gene expression in human melanoma cells⁴⁸, and with the notion that the methylation of some histone
402 residues depends on the transcription machinery⁴³. While we observed that certain modifications, such
403 as H3K4me1/2 and H3K27ac tend to anticipate gene expression, this does not necessarily mean that they
404 are the cause of transcription initiation. Actually, we have also observed particular cases in which these
405 marks are deposited post-activation (for an example see Figure 5b, middle panels). After the initial stage
406 of gene activation, further changes in gene expression, even if substantial, appear to be mostly uncoupled
407 from changes in histone modifications (Figure 7b). It is tempting to speculate that after the initial burst of
408 transcription, histone residues are saturated with modifications, and that therefore, any further up-regulation
409 of gene expression cannot possibly be accompanied by increased levels of histone modifications.

410 We do have identified a small set of genes that are expressed in the absence of any histone modification,
411 with the exception of H3K4me1 and H3K4me2 (Figures 4c, 5a-b lower panels). A few of these are activated
412 later during the transdifferentiation process, and therefore we lack the temporal resolution to detect post-
413 activation marking. Still, many of these genes are down-regulated or stably expressed, and are unmarked
414 even at the beginning of transdifferentiation (for an example see Figure 5b, lower panels). Gene activation
415 without histone modifications has been previously observed for developmentally regulated genes in the fruit
416 fly¹⁵.

417 Here we have focused specifically on the dynamics of chromatin modifications during up-regulation.
418 Our results suggest that down-regulation appears to be largely uncoupled from chromatin changes (Sup-
419 plementary Figure 5h). However, while RNA sequencing-inferred expression levels can be used to approx-
420 imately identify the time at which a gene is initially activated, differences in RNA stability may confound
421 the identification of the time-point at which a gene is fully inactivated. Indeed, RNAs can be detected long
422 after gene inactivation, for a time likely to be specific to each individual gene. Therefore, the data that
423 we have generated does not have the appropriate resolution to discard that this lack of coupling during
424 down-regulation is partially caused by the difficulty in precisely identifying the time-point at which genes
425 stop being expressed.

426 The multi-omics data that we have generated during the pre-B cell transdifferentiation into macrophages
427 has allowed us to address with unprecedented resolution some fundamental questions regarding the dy-
428 namics of chromatin marking and gene expression during cellular differentiation, and have contributed to
429 shed light on some long-standing questions in the field. These findings may have implications on therapeu-
430 tic strategies currently relying on the causal role of chromatin modifications⁹⁻¹¹. Further mining of this data
431 resource will certainly contribute to a deeper understanding of the epigenetic layer of gene regulation.

432 **Methods**

433 **RESOURCE AVAILABILITY**

434 **Materials Availability**

435 This study did not generate new unique reagents.

436 **Data and Code Availability**

437 The code generated during this study is available at [https://github.com/bborsari/Borsari_et_](https://github.com/bborsari/Borsari_et_al_transdifferentiation_chromatin)
438 [al_transdifferentiation_chromatin](https://github.com/bborsari/Borsari_et_al_transdifferentiation_chromatin). A complete list of scripts used for each analysis described
439 in the section *Method details* can be found at [https://github.com/bborsari/Borsari_et_al_](https://github.com/bborsari/Borsari_et_al_transdifferentiation_chromatin/blob/master/bin/table.scripts.tsv)
440 [transdifferentiation_chromatin/blob/master/bin/table.scripts.tsv](https://github.com/bborsari/Borsari_et_al_transdifferentiation_chromatin/blob/master/bin/table.scripts.tsv). When not speci-
441 fied in the text, the code used for a given analysis is included in the corresponding figure's script.

442 RNA-seq and ChIP-seq raw and processed data from this study have been submitted to ArrayExpress

443 (<https://www.ebi.ac.uk/arrayexpress/>) under accession numbers E-MTAB-9790 and E-MTAB-
444 9825, respectively.

445 Processed data in GRCh38/hg38 assembly from this study is available for visualization at the UCSC
446 Genome Browser⁴⁹ (<http://genome.ucsc.edu/>). The track data hub is available at
447 https://public-docs.crg.es/rquigo/Data/bborsari/hubs/ERC_human_hub/hub.txt.

448 A web page has also been implemented to gather all information regarding the Chromatin and Tran-
449 scriptomics Dynamics Project (<http://rnamaps.crg.eu/>). The web page provides information about
450 all experiments and replicates performed during the project, as well as access to the data in ArrayExpress
451 and the UCSC Genome Browser.

452 ENCODE data is freely available on the ENCODE portal (<https://www.encodeproject.org/>).
453 Experiments and files accession IDs for RNA-seq and ChIP-seq data are reported in Supplementary Tables
454 5 and 6, respectively.

455 **EXPERIMENTAL MODEL AND SUBJECT DETAILS**

456 **Transdifferentiation of BLAER1 cells to macrophages**

457 For the transdifferentiation process we made use of the Burkitt lymphoma cell line BlaER1, as described
458 in 24. Induction of transdifferentiation (treatment with 100 μ M β -estradiol and growth in the presence of 10
459 nM Il-3 and 10 nM CSF-1) has been described in 50 and 51. The process was monitored at 12 time-points
460 (as described in 24): 0, 3, 6, 9, 12, 18, 24, 36, 48, 72, 120 and 168 hours post-induction (p.i.; Figure 1a).

461 **METHOD DETAILS**

462 **RNA-seq library preparation and sequencing**

463 Two independent biological replicates for each time-point were performed. Briefly, cells were lysed with
464 QiAZol (Qiagen, The Netherlands). Chloroform was added to each sample, and RNA contained in the
465 aqueous solution was isolated and purified by using RNeasy mini kit columns (Qiagen, The Netherlands).
466 Poly A+ libraries were prepared with 1 μ g of total RNA and using TruSeq Stranded mRNA Library Prep
467 Kit (Illumina, USA) according to the manufacturer's protocol. Libraries were analyzed using Agilent DNA
468 1000 chips to determine the quantity and size distribution, and sequenced paired-end 75-bp on an Illumina
469 HiSeq 2000.

470 **ChIP-seq library preparation and sequencing**

471 ChIP-seq experiments of nine histone marks (H3K4me1: Abcam ab8895; H3K4me2 : Millipore 07-030;
472 H3K4me3: Abcam ab8580; H3K9ac: Abcam ab4441; H3K27ac: Diagenode C15410192; H3K36me3:
473 Abcam ab9050; H4K20me1: Abcam ab9051; H3K9me3: Abcam ab8898; H3K27me3: Millipore 07-449)
474 were performed in two independent biological replicates for each time-point. Cells were crosslinked with

475 formaldehyde 1% (Sigma) for 10' at room temperature. The reaction was stopped by adding glycine to
476 0.25 M final concentration for 10' at room temperature. Fixed cells were resuspended in 100 μ L of lysis
477 buffer (SDS 1%, EDTA 10 mM, TrisCl 50 mM and protease inhibitors). The lysate was sonicated for 25'
478 using Covaris S2 system in TC12 tubes (Duty cycle 20%, Intensity 8, cycles/burst 200, water level 15).
479 The cleared supernatant was used immediately in ChIP experiments or stored at -80 °C. 5 μ g of sonicated
480 chromatin were diluted in 900 μ L RIPA buffer — H3K4me3, H3K9ac, H4K20me1, H3K27me3 and H3K27ac
481 (140 mM NaCl, 10 mM Tris-HCl pH 8.0, 1 mM EDTA, 1% Triton X-100, 0.1% SDS, 0.1% Na deoxycholate,
482 protease inhibitors) —, RIPA 2X — H3K4me1, H3K4me2 and H3K9me3 (280 mM NaCl, 10 mM Tris-
483 HCl pH 8.0, 1 mM EDTA, 2% Triton X-100, 0.2% SDS, 0.2% Na deoxycholate, protease inhibitors) —,
484 or RIPA 1X 1% triton — H3K36me3 (280 mM NaCl, 10 mM Tris-HCl pH 8.0, 1 mM EDTA, 1% Triton X-
485 100, 0.2% SDS, 0.2% Na deoxycholate, protease inhibitors). For H3K4me3, H3K36me3, H3K9ac and
486 H3K27me3 ChIPs, chromatin and antibodies were incubated overnight, rotating at 4 °C with 0.125-5 μ g of
487 specific antibody and samples were then incubated for 2 hours rotating at 4 °C with Dynabeads protein A
488 for immunoprecipitation (Invitrogen) to recover the bound material. For H3K4me1, H3K4me2, H3K9me3,
489 H4K20me1 and H3K27ac ChIPs, antibodies were coated to protein A magnetic beads for 2 hours at 4
490 °C prior to overnight incubation with chromatin. In all cases, beads were washed for 10' three times in 1
491 mL of the corresponding immunoprecipitation buffer without protease inhibitors, then washed once in 1 mL
492 LiCl buffer (0.25 M LiCl, 0.5% NP-40, 0.5% sodium deoxycholate, 1 mM Na-EDTA, 10 mM Tris-HCl, pH
493 8.0), and finally washed twice in 1 mL of TE buffer (1 mM Na-EDTA, 10 mM Tris-HCl, pH 8.0). ChIPped
494 material was incubated with DNase-free RNase at 50 μ g/mL for 30' at 37 °C. Chromatin was reverse-
495 crosslinked by adding SDS (0.5% final concentration) and Proteinase K (500 μ g/mL final concentration)
496 and incubated overnight at 65 °C. ChIPped chromatin was then purified with Qiaquick PCR purification
497 columns (Qiagen) following the manufacturer's instructions. ChIP libraries were prepared with 1-5 ng of
498 DNA and using NebNext Ultra DNA library prep kit for Illumina (New England Biolabs) according to the
499 manufacturer's protocol. Libraries were analyzed using Agilent DNA High Sensitivity chips to determine the
500 quantity and size distribution, and sequenced single-read 50-bp on an Illumina HiSeq 2000.

501 In total, 264 samples were sequenced (24 by RNA-seq, 216 by ChIP-seq, 24 by ChIP input).

502 RNA-seq data processing and analysis

503 Data was processed using the *grape-nf* (<https://github.com/guigolab/grape-nf>) Nextflow²⁵
504 pipeline. RNA-seq reads were aligned to the human genome (assembly GRCh38, Gencode annotation
505 version 24) using the STAR⁵² software version 2.4.0j . We allowed a maximum number of mismatches
506 equal to 4% of the read length. Only alignments for reads mapping to ten or fewer loci were reported.
507 Quantification of genes and transcripts was done with RSEM⁵³ version 1.2.21. TPM calculation was per-
508 formed after removing mitochondrial genes.

509 From the set of 19,831 protein-coding genes (Gencode v24), we selected 10,696 expressed genes with
510 a maximum expression during transdifferentiation \geq 5 TPM in both replicates, and 1,552 silent genes (0

511 TPM in all time-points and replicates). Based on this set of 12,248 genes, we quantile-normalized the ex-
512 pression matrices (\log_2 -transformed TPM, pseudocount of 1) across replicates and time-points using the R
513 package `preprocessCore`⁵⁴ (script: `quantile.normalization.R`), and obtained the mean expression
514 levels between replicates (script: `matrix.matrix.mean.R`).

515 To detect significant gene expression changes along transdifferentiation, we used the R package `maSig-`
516 `Pro`⁵⁵ with replicates handled internally. Function `p.vector()` was run with default parameters: $Q = 0.05$,
517 `MT.adjust = "BH"`, `min.obs = 20` (script: `maSigPro.wrapper.R`). We defined as stably expressed
518 those genes reporting a `maSigPro` FDR value ≥ 0.05 ($n = 2,666$).

519 As concerns the identification of up-regulated, down-regulated, peaking and bending genes, we per-
520 formed a two-step classification across the 8,030 genes with significantly variable gene expression profiles.
521 Briefly, we first focused on profiles with at least two-fold change (in \log_2 scale this change corresponds to 1)
522 and identified monotonic up-regulations and down-regulations; peaking profiles were defined as monotonic
523 increases followed by monotonic decreases, bending profiles as the opposite (script: `classification.`
524 `log2.pl`). All other significantly variable genes with fold-change < 2 were assigned to one of these four
525 groups following hierarchical clustering (distance measure: *euclidean*; clustering method: *complete*; script:
526 `classification.2.R`).

527 **ChIP-seq data processing and analysis**

528 Data was processed using the *ChIP-nf* (<https://github.com/guigolab/chip-nf>) Nextflow²⁵
529 pipeline. ChIP-seq reads were aligned to the human genome assembly (GRCh38) using the GEM⁵⁶ map-
530 ping software, allowing up to two mismatches. Only alignments for reads mapping to ten or fewer loci
531 were reported. Duplicated reads were removed using Picard ([http://broadinstitute.github.io/
532 picard/](http://broadinstitute.github.io/picard/)). Pile-up signal from bigWig files was obtained running MACS2⁵⁷ on individual replicates. No
533 shifting model was built. Instead, fragment length was set to 250 bp and was used to extend each read
534 towards the 3' end (using the `--extsize` option). Pile-up signal was normalized by scaling larger sam-
535 ples to smaller samples (using the default for the `--scale-to` option) and adjusting signal per million
536 reads (enabling the `--SPMR` option). Peak calling was performed using Zerone⁵⁸ with replicates handled
537 internally, and passed the filter for all pairs of replicates (advice: `accept discretization`).

538 To check library complexity, we computed the fraction of non-redundant mapped reads⁵⁹ (recom-
539 mended threshold: $NRF \geq 0.8$) for each ChIP-seq experiment, and found a minimum NRF value of 0.92.
540 Additionally, to evaluate the global ChIP enrichment, we computed the fraction of reads in peaks⁵⁹ (recom-
541 mended threshold: $FRiP \geq 0.01$), and found a minimum FRiP value of 0.05.

542 The intersection / overlap analyses described below were performed with the function `intersectBed`
543 of BEDTools⁶⁰ software v2.27.1.

544 To select the genomic location enriched, on average, in a specific histone mark (region of interest),
545 we focused on an up-stream and down-stream 5 Kb region (± 5 Kb) with respect to the first annotated
546 Transcription Start Site (TSS) of the gene, and retrieved 6,063 protein-coding genes that did not overlap

547 any other gene body ± 5 Kb. For each histone modification we then selected, among the 6,063 genes,
548 those with peaks in the ± 5 Kb promoter region in all the 12 time-points, and computed, using the function
549 `aggregate` from the `bwtool`⁶¹ software (script: `bwtool.aggregate.ChIPseq.sh`), the mean pile-up
550 signal for each experiment. Based on this analysis, we decided to select as regions of interest i) the gene
551 body for H3K36me3 and H4K20me1, ii) ± 2 Kb with respect to the TSS for all other marks (Supplementary
552 Figure 1c). A comprehensive catalogue of all non-redundant (same ensembl gene ID and start coordinate)
553 TSSs annotated for the selected 12,248 in Gencode v24 was obtained with the script `non.redundant.`
554 `TSS.sh`.

555 To compare expression and chromatin profiles over time, we quantified, for each of the nine histone
556 marks, the amount of pile-up signal associated with a gene at each time-point (script: `get.matrix.`
557 `chipseq.sh`). Briefly, if a peak was present in the region of interest of a gene at a specific time-point,
558 we considered the mean pile-up signal in the intersection between the peak and the region of interest,
559 otherwise we computed the mean pile-up value in the entire region of interest. In the presence of multiple
560 peaks and/or multiple regions of interest (e.g. in case of multiple TSSs annotated for the same gene), we
561 considered the highest of all observed values. Matrices of histone marks' signals for the selected 12,248
562 protein-coding genes were quantile-normalized across replicates and time-points using the R package
563 `preprocessCore`⁵⁴ as done for gene expression. For all down-stream analyses, we used the mean signal
564 between replicates.

565 **Principal Component Analysis of expression and chromatin data**

566 For this type of analysis we made use of the transposed expression and chromatin For this type of analysis
567 we made use of the transposed expression and chromatin matrices generated as described in sections
568 *RNA-seq data processing and analysis* and *ChIP-seq data processing and analysis*, respectively. There-
569 fore, genes (columns) and time-points (rows) were used as variables and observations, respectively. We
570 centered and scaled each of the ten transposed matrices independently, obtaining z-score profiles for each
571 time-point monitored at expression and histone marks' level. For the joint Principal Component Analysis
572 (PCA) reported in Figure 1c across expression and the nine histone marks, we included as variables the
573 subset of 10,658 genes with non-missing (NA) z-score profiles in all ten matrices. As a consequence,
574 1,590 genes were excluded from this analysis, 98% of them being the silent genes (1,552). For the PCAs
575 reported in Supplementary Figure 3d, we considered for each histone modification the corresponding sets
576 of DE genes that are either stably or differentially marked.

577 **Analysis of the degree of correlation between expression levels and chromatin signals**

578 Steady-state correlations between gene expression levels and each histone mark's signals were computed
579 at individual time-points considering the entire set of 12,248 selected protein-coding genes. In this case,
580 Pearson r measured the degree of correlation between the vector of 12,248 expression levels and the vector
581 of 12,248 mark signals at a given time-point (Figure 1d, dots). Time-course correlations were measured,

582 instead, at the level of individual expressed genes. Silent genes were not considered for this analysis,
583 because of the zero standard deviation in their time-series expression profile (i.e. 0 TPM in all time-points).
584 Thus, for each gene and histone mark we obtained the Pearson r correlation coefficient between the vector
585 of 12 expression levels (i.e. the expression levels measured at the 12 time-points) and the vector of 12 mark
586 signals. The distributions of Pearson r correlation coefficients for the set of (differentially + stably) expressed
587 genes are depicted with box plots and violin plots in Figure 1d. Randomized steady-state and time-course
588 correlation coefficients were computed as described above following a 1,000-permutations scheme on each
589 histone mark's matrix. Briefly, while we kept the original expression matrix, the columns (time-points) of
590 the matrix corresponding to a given mark's signal were permuted without repetition 1,000 times (for an
591 example, see Supplementary Figure 2a, lower panel). In the case of steady-state correlations we report,
592 for each expression time-point, the Pearson r averaged over 1,000 rounds of permutation of chromatin
593 time-points (Supplementary Figure 2b, dots). In the case of correlations computed across time-points (time-
594 course), we computed, for each gene, the Pearson r averaged over the 1,000 rounds of permutations. The
595 distributions of the resulting coefficients across the set of expressed genes are depicted in Supplementary
596 Figure 2b (box plots and violin plots). Correlations were computed with the R function `cor()`. Permutations
597 without replacement of the chromatin time-points were performed consistently across histone marks with
598 the R function `sample()`, by setting an independent seed for each round of permutations. The correlation
599 values reported in Supplementary Figure 2c are an analogous exercise to Figure 1d on the set of 8,030
600 differentially expressed genes.

601 **Multivariate Hidden Markov Model analysis**

602 A multivariate Hidden Markov Model (HMM) was fitted to the entire ChIP-seq dataset to approximate the
603 set of underlying chromatin states reported by the 12,248 selected protein-coding genes along the transd-
604 ifferentiation process. Specifically, we provided as input a matrix of dimensions 146,976 rows \times 9 columns,
605 which collected for each gene and time-point (12,248 genes, 12 time-points) the signal of each of the
606 9 histone marks after quantile normalization (for a description of these calculations see previous section
607 *ChIP-seq data processing and analysis*). The collective behavior of the nine histone marks along the twelve
608 time-points was modelled as an independent time-series for each gene, using Gaussian distributions. The
609 model then reprocessed each gene's data to estimate the chromatin state of each gene at each time-point,
610 and provide a time series of chromatin states for each gene. HMM was performed using the R package `dep-`
611 `mixS4`⁶², in particular functions `depmix()`, `fit()` and `posterior()` (script: `HMM.wrapper.marks.R`).
612 We repeated the analysis for increasing numbers of states (between 2 and 20), and recorded the log
613 likelihood of each model (the 20-states model reached the maximum number of iterations in EM without
614 convergence). We found that somewhere between five and eight states approximate the elbow point of
615 the log likelihood curve (Supplementary Figure 3a), and observed that the combinations of histone marks
616 represented by five states were consistent with manual inspection of pile-up histone marks profiles in the
617 UCSC genome browser. We thus set for five states. The response parameters of the nine histone marks

618 corresponding to each of these states are reported in Figure 2a. In this case, the *Intercept* values of each
619 histone mark across the five states were re-scaled to a range 0-1 to enable the comparison among differ-
620 ent states and marks. HMM sequence hierarchical clustering across the 12,248 genes was performed with
621 the TraMineR⁶³ and pheatmap (<https://github.com/raivokolde/pheatmap>) R packages (cluster-
622 ing distance: *euclidean*, clustering method: *Ward.D2*). The arc diagram representation in Figure 2c was
623 obtained with the R package arcdiagram (<https://github.com/gastonstat/arcdiagram>).

624 **Decision-tree labelling**

625 In the Methods section *ChIP-seq data processing and analysis* we introduced the distinction between genes
626 with and without peaks of a given mark at a given point in the region of interest (gene body for H3K36me3
627 and H4K20me1; TSS ± 2 Kb for all other marks). Following this first assessment, we classified as unmarked
628 those genes that were consistently unmarked throughout the whole process of transdifferentiation, i.e. with
629 no peaks called at any time-point in the region of interest. Conversely, marked genes reported peak calls
630 of a given mark in the region of interest in at least one time-point (Figure 4a).

631 Within the set of marked genes, we defined as stably marked (SM) those that did not report significant
632 changes detected by maSigPro⁵⁵ over time ($FDR \geq 0.05$). On the contrary, differentially marked (DM)
633 genes reported significant changes in a given mark's profile over time ($FDR < 0.05$). To ensure a multiple
634 testing correction procedure consistent among the nine marks and also with respect to gene expression,
635 maSigPro was run, as described for gene expression (default parameters, replicates handled internally),
636 on the initial set of 12,248 genes, which also included unmarked genes.

637 The next branch of classification (Figure 4a) was applied only to the set of differentially marked genes
638 that are also differentially expressed. To ensure consistent results among histone marks, the following mul-
639 tiple testing correction procedures were always applied to the set of 8,030 DE genes. For each DE gene,
640 we computed at each time-point the breadth of a given mark's signal, defined as the fraction of the gene's
641 size (from the first annotated region of interest until the last annotated Transcription Termination Site, TTS)
642 covered by peaks of the mark. We refer to this vector of length 12 as the mark's coverage vector. We
643 next considered i) Pearson *r* correlation coefficient between the time-series expression levels and mark's
644 signals; ii) Pearson *r* correlation coefficient between the time-series expression levels and mark's coverage
645 values; iii) statistical significance of the Needleman-Wunch (NW) dynamic time warping alignment be-
646 tween the time-series expression levels and mark's signals (following Benjamini-Hochberg multiple testing
647 correction; script: `p-adjust.R`). We used as input for the NW alignments (scripts: `NW.alignment.`
648 `path.R`, `NW.bidirectional.matches.py`) the z-score profiles of expression and mark obtained after
649 applying polynomial regression (`degree = 2`) on the original matrices (scripts: `loess.polynomial.`
650 `regression.R`, `NW.generate.input.matrix.sh`). This procedure was applied to remove the noise
651 due to occasional fluctuations in signal over time. A permutation *p* value for each gene was computed
652 (script: `NW.pvalue.permutation.test.py`), based on a 100,000-permutations scheme (script: `NW.`
653 `alignment.permutations.R`). To classify a gene as positively correlated, we required at least two of

654 the following conditions: i) Pearson r correlation coefficient between the time-series expression levels and
655 mark's signals ≥ 0.60 and $FDR < 0.05$; ii) Pearson r correlation coefficient between the time-series expres-
656 sion levels and mark's coverage values ≥ 0.60 and $FDR < 0.05$; iii) NW alignment between the time-series
657 expression levels and mark's signals with $FDR < 0.05$. For negatively correlated genes, we required at
658 least two of the following conditions: i) Pearson r correlation coefficient between the time-series expres-
659 sion levels and mark's signals ≤ -0.60 and $FDR < 0.05$; ii) Pearson r correlation coefficient between the
660 time-series expression levels and mark's coverage values ≤ -0.60 and $FDR < 0.05$; iii) NW alignment be-
661 tween the time-series expression levels and mark's signals with $FDR \geq 0.05$. Genes that did not meet
662 these requirements were classified as uncorrelated. The same decision-tree classification was performed
663 independently for each of the nine histone marks, to ensure comparable results among all modifications
664 (`script: define.6.groups.R`).

665 **Clustering analysis**

666 We considered all 45 combinations between the 9 histone marks and the 5 decision-tree labels described
667 in the previous section. For instance, one combination may be “stably marked + H3K4me3”, and another
668 combination may be “positively correlated + H3K27ac”. To test the co-occurrence of this pair of combi-
669 nations, we retrieved the set of DE genes that are labelled “stably marked” for H3K4me3, and the set of
670 DE genes that are labelled “positively correlated” for H3K27ac. The significant overlap between these two
671 sets of genes was tested by the hypergeometric distribution (R function `phyper()`). We repeated this
672 procedure for all possible pairs of combinations. We next clustered the p values obtained after applying the
673 Benjamini-Hochberg False Discovery Rate (FDR) multiple testing correction. Hierarchical clustering was
674 performed with the `ComplexHeatmap`⁶⁴ R package (clustering distance = *Manhattan*, clustering method
675 = *Ward.D2*). Cluster correspondence analysis⁶⁵ of the 45 categorical variables (combinations of histone
676 marks and decision-tree labels) across the 8,030 selected genes was performed with the R package `clus-`
677 `trd`⁶⁶. To select the optimal number of clusters and dimensions, we first run the function `tuneclus()`
678 with the following parameters: `nclusrange = 3:10`, `ndimrange = 2:9`, `method = "clusCA"`, `nstart`
679 `= 100`, `seed = 1234`. This indicated that the optimal number of dimensions and clusters was two and
680 three, respectively. We then obtained the three clusters of genes running the function `clusmca` with the
681 following parameters: `nclus = 3`, `ndim = 2`, `method = "clusCA"`, `nstart = 100`, `smartStart = NULL`,
682 `gamma = TRUE`, `seed = 1234`. We obtained the same clusters of genes when running the function `clusmca`
683 with the following parameters: `nclus = 3`, `ndim = 3`, `method = "MCAk"`, `alphak = 0.5`, `nstart = 100`,
684 `smartStart = NULL`, `gamma = TRUE`, `seed = 1234`). This allowed us to explore the clustering of genes
685 also in the third dimension (Figure 4c, Supplementary Figure 4a).

686 **Gene Ontology enrichment analysis**

687 We used the R package `GOstats`⁶⁷ to identify Gene Ontology (GO) terms related to biological processes
688 (BP) and cellular compartments (CC). We set a p value threshold of 0.01 to identify significantly enriched

689 terms. For the GO enrichment analysis on the genes contributing to Principal Components (PC) 1 and
690 2 (described in Results, section *Gene expression recapitulates transdifferentiation more precisely than*
691 *chromatin*; Figure 1c, Supplementary Table 2), we used the function `get_pca_var()` from the R pack-
692 age `factoextra` (<https://CRAN.R-project.org/package=factoextra>) to extract the 10% genes
693 ($n = 1,066$) with the highest contribution to each of the two first principal components. The union of
694 these two sets of genes was used as background for the GO enrichment analysis. We used REVIGO⁶⁸
695 (<http://revigo.irb.hr/>) to summarize the lists of enriched GO terms. For the GO enrichment anal-
696 ysis on the up-regulated genes that belong to the three chromatin clusters (described in Results, section
697 *Chromatin marking is associated with expression specifically at the time of gene activation*), we provided
698 as background the set of 2,103 up-regulated genes. In this case, we used REVIGO and the R package
699 `ggplot2`⁶⁹ to compute and visualize, respectively, maps of the identified GO terms based on their frequency,
700 $-\log_{10} p$ value, uniqueness and dispensability. Only children terms with dispensability < 0.5 are shown.

701 **Analysis of ENCODE RNA-seq and ChIP-seq data**

702 To investigate differences in gene expression levels and chromatin marking among the three clusters of
703 DE genes in other biological models, we obtained RNA-seq data and ChIP-seq data for histone marks
704 generated by the ENCODE Project^{70,71} (<https://www.encodeproject.org/>). Besides B cells and
705 CD14-positive monocytes, which are biologically more similar to pre-B cells and macrophages, respec-
706 tively, we selected five cancer cell lines (K562, HepG2, GM12878, MCF-7, A549) that are comprehensively
707 characterized by ENCODE ChIP-seq data for the nine histone marks that we have profiled in our study. To
708 assess differences in gene expression levels between the three clusters of DE genes, we obtained gene ex-
709 pression quantifications (with respect to Gencode v24) from polyA+ RNA-seq experiments (accession date:
710 10/06/2019). We computed, for each gene, the average TPM values between two biological replicates. The
711 list of experiments and datasets' accession IDs used for this analysis is reported in Supplementary Table
712 5.

713 To assess differences in chromatin marking, we obtained ChIP-seq data available for the nine histone
714 marks profiled in our study. (Assay title: Histone ChIP-seq; Genome assembly: GRCh38; Output type:
715 replicated peaks or stable peaks; Accession date: 10/06/2019). The list of experiments and datasets'
716 accession IDs used for this analysis is available in Supplementary Table 6. In all cases, we excluded
717 experiments associated with AUDIT errors. In case of multiple experiments on the same target and cell
718 type, the experiment associated with the lowest number of AUDIT terms was selected. The scripts used to
719 retrieve and filter the ENCODE experiments are: `download.metadata.sh`, `parse.metadata.audit`,
720 `categories.py`, `retrieve.encode.identifiers.sh`, `parse.list.identifiers.sh`.

721 For each experiment and cell type, we computed the proportion of genes with at least one peak called
722 over the gene body (H3K36me3, H4K20me1) or in the promoter region (TSS ± 2 Kb for all other marks;
723 script: `intersect.peaks.regions.sh`). In the presence of multiple TSSs annotated for the same
724 gene, multiple regions were considered. This is consistent with the analyses described in section *ChIP-seq*

726 **Analysis of temporal dynamics**

727 For this analysis we first identified, within the set of 2,103 up-regulated genes, 257 with expression at 0
728 hours p.i. < 1 TPM. These genes were, therefore, specifically activated during transdifferentiation. Ex-
729 pression and chromatin profiles of each of the considered genes were re-scaled to range 0-100 (script:
730 `rescale.R`): in this way, the minimum and maximum expression level or chromatin signal over the 12
731 time-points were set to 0% and 100% of up-regulation, respectively. We next considered, for each gene,
732 pairs of consecutive time-points along transdifferentiation (e.g. 0h and 3h; 3h and 6h; 6h and 9h; etc.),
733 and recorded the first time-point at which the expression / chromatin profile crossed (\geq) 25%, 50%, 75%
734 and 100% degree of up-regulation (Supplementary Figure 5b). This “crossing” step implies that, in a pair
735 of consecutive time-points, the signal corresponding to the first time-point is, for instance, $< 25\%$, and the
736 signal corresponding to the second time-point is, for instance, $\geq 25\%$. This assessment is performed for
737 each of the four degrees of up-regulation. To ensure monotonic increases consistently across all histone
738 marks, we excluded genes for which this “crossing” step could not be observed for all four degrees of
739 up-regulation in a given mark’s time-series profile. This explains the different numbers of genes, among
740 marks, reported in Figure 6a and Supplementary Figure 5e. For a given gene and for each of the four de-
741 grees of up-regulation, the recorded time-points (tp) for expression and chromatin profiles were compared,
742 and a label was assigned depending on whether the up-regulation of chromatin signal anticipated (tp_{mark}
743 $< tp_{expression}$), co-occurred ($tp_{mark} = tp_{expression}$) or followed ($tp_{mark} > tp_{expression}$) the up-regulation of
744 gene expression. We analogously compared the up-regulation between pairs of histone marks (Figure 6c,
745 Supplementary Figure 5d). In this case, we analyzed whether the up-regulation of histone mark’s signal
746 on row i anticipated ($tp_i < tp_j$) or co-occurred with ($tp_i = tp_j$) the up-regulation of histone mark’s signal on
747 column j . To assess whether the specific order of up-regulation in expression levels and chromatin signals
748 depended on the initial level of expression of the genes, these analyses were repeated starting on a set of
749 629 up-regulated genes with expression at 0 hours p.i. > 25 TPM.

750 **QUANTIFICATION AND STATISTICAL ANALYSIS**

751 Details regarding statistical tests, significance assessment, dispersion and precision measures are re-
752 ported both in the section *Method details* and in the figures’ legends. All statistical analyses were performed
753 using the R language for statistical computation and graphics⁷²(<http://www.R-project.org/>). In all
754 cases, the multiple testing correction procedure was performed by applying the Benjamini-Hochberg⁷³
755 False Discovery Rate (FDR). Wilcoxon rank-sum tests were performed with the `wilcox.test()` R func-
756 tion in a two-sided manner.

757 When not specified, plots were made using the R package `ggplot2`⁶⁹. All box plots depict the first and
758 third quartiles as the lower and upper bounds of the box, with a band inside the box showing the median

759 value and whiskers representing 1.5x the interquartile range. All scripts used in the analyses are publicly
760 available (see the *Data and Code Availability* statement).

761 **Acknowledgments**

762 We thank Thomas Graf and Francesca Rapino for donating BLaER1 cells and for helpful discussions. We
763 thank Sebastian Ullrich, Carme Arnan and Vasilis Ntasis for helpful discussion about the data. We thank
764 Montserrat Corominas, Guillaume Fillion and Luciano Di Croce for insightful suggestions. We thank Diego
765 Garrido-Martín, Manuel Muñoz and Javier Martín-Vallejo for statistical advice. We also thank the Genomics,
766 the Flow Cytometry and the Bioinformatics Core Units of the CRG (Barcelona, Spain). We thank Romina
767 Garrido for administrative support. We thank the ENCODE Consortium, in particular Thomas Gingeras',
768 Bradley Bernstein's, John Stamatoyannopoulos' and Peggy Farhnam's laboratories, for data production.
769 This work was performed under the financial support of the European Community under the FP7 pro-
770 gram (ERC-2011-AdG-294653-RNA-MAPS). B.B. is supported by the fellowship 2017FI_B00722 from the
771 Secretaria d'Universitats i Recerca del Departament d'Empresa i Coneixement (Generalitat de Catalunya)
772 and the European Social Fund (ESF). C.C.K. is supported by the CERCA Programme / Generalitat de
773 Catalunya and FEDER under project VEIS-001-P-001647. B.R.C. is supported by the Ministerio de Cien-
774 cia, Innovación y Universidades de España under grant FJCI-2017-34353. We also acknowledge Agencia
775 Estatal de Investigación (AEI) and FEDER under project PGC2018-094017-B-I00. All authors acknowledge
776 the support of the Ministerio de Ciencia, Innovación y Universidades de España to the EMBL partnership,
777 the Centro de Excelencia Severo Ochoa, and the CERCA Programme / Generalitat de Catalunya. Figures
778 1b and 7a were created with <https://biorender.com>.

779 **Author Contributions**

780 R.G. and R.J. conceived the project. B.B., S.P-L. and R.G. designed the study. B.B. performed the com-
781 putational analyses. A.A. performed the ChIP-seq experiments. A.E. and M.S. performed the RNA-seq
782 experiments. C.C.K. and E.P. contributed to data quality check and processing. C.C.K., R.N., M.R-R. and
783 B.R.C. contributed tools and ideas to perform experiments and computational analyses. B.B., S.P-L. and
784 R.G. wrote the manuscript with the contribution of all authors.

785 **Competing Interests**

786 The authors declare no competing interest.

787 **References**

- 788 [1] Hon, G., Wang, W. & Ren, B. Discovery and Annotation of Functional Chromatin Signatures in the
789 Human Genome. *PLoS Computational Biology* **5**, e1000566 (2009).
- 790 [2] Barski, A. *et al.* High-Resolution Profiling of Histone Methylations in the Human Genome. *Cell* **129**,
791 823–837 (2007).
- 792 [3] Schneider, R. *et al.* Histone H3 lysine 4 methylation patterns in higher eukaryotic genes. *Nature Cell*
793 *Biology* **6**, 73–77 (2004).
- 794 [4] Trojer, P. & Reinberg, D. Facultative Heterochromatin: Is There a Distinctive Molecular Signature?
795 *Molecular Cell* **28**, 1–13 (2007).
- 796 [5] Hansen, K. H. *et al.* A model for transmission of the H3K27me3 epigenetic mark. *Nature Cell Biology*
797 **10**, 1291–1300 (2008).
- 798 [6] Strahl, B. D. & Allis, C. D. The language of covalent histone modifications. *Nature* **403**, 41–45 (2000).
- 799 [7] Karlič, R., Chung, H. R., Lasserre, J., Vlahoviček, K. & Vingron, M. Histone modification levels are
800 predictive for gene expression. *Proceedings of the National Academy of Sciences of the United States*
801 *of America* **107**, 2926–2931 (2010).
- 802 [8] Dong, X. *et al.* Modeling gene expression using chromatin features in various cellular contexts.
803 *Genome Biology* **13**, R53 (2012).
- 804 [9] Guo, S. *et al.* Epigenetic Regulation Mediated by Methylation in the Pathogenesis and Precision
805 Medicine of Rheumatoid Arthritis. *Frontiers in Genetics* **11**, 811 (2020).
- 806 [10] Rahman, M. M. & Tollefsbol, T. O. Targeting cancer epigenetics with CRISPR-dCAS9: Principles and
807 prospects. *Methods* (2020).
- 808 [11] Zhao, W., Wang, Y. & Liang, F. S. Chemical and light inducible epigenome editing. *International*
809 *Journal of Molecular Sciences* **21**, 998 (2020).
- 810 [12] Greer, E. L. & Shi, Y. Histone methylation: A dynamic mark in health, disease and inheritance. *Nature*
811 *Reviews Genetics* **13**, 343–357 (2012).
- 812 [13] Nègre, N. *et al.* A cis-regulatory map of the Drosophila genome. *Nature* **471**, 527–531 (2011).
- 813 [14] Hödl, M. & Basler, K. Transcription in the absence of histone H3.2 and H3K4 methylation. *Current*
814 *Biology* **22**, 2253–2257 (2012).
- 815 [15] Pérez-Lluch, S. *et al.* Absence of canonical marks of active chromatin in developmentally regulated
816 genes. *Nature Genetics* **47**, 1158–1167 (2015).

- 817 [16] Vandenbon, A., Kumagai, Y., Lin, M., Suzuki, Y. & Nakai, K. Waves of chromatin modifications in
818 mouse dendritic cells in response to LPS stimulation. *Genome Biology* **19**, 138 (2018).
- 819 [17] Le Martelot, G. *et al.* Genome-Wide RNA Polymerase II Profiles and RNA Accumulation Reveal
820 Kinetics of Transcription and Associated Epigenetic Changes During Diurnal Cycles. *PLoS Biology*
821 **10**, e1001442 (2012).
- 822 [18] Wang, S. *et al.* A dynamic and integrated epigenetic program at distal regions orchestrates transcrip-
823 tional responses to VEGFA. *Genome Research* **29**, 193–207 (2019).
- 824 [19] Rach, E. A. *et al.* Transcription Initiation Patterns Indicate Divergent Strategies for Gene Regulation
825 at the Chromatin Level. *PLoS Genetics* **7**, e1001274 (2011).
- 826 [20] Mercer, E. M. *et al.* Multilineage Priming of Enhancer Repertoires Precedes Commitment to the B and
827 Myeloid Cell Lineages in Hematopoietic Progenitors. *Immunity* **35**, 413–425 (2011).
- 828 [21] Kaikkonen, M. U. *et al.* Remodeling of the enhancer landscape during macrophage activation is
829 coupled to enhancer transcription. *Molecular Cell* **51**, 310–325 (2013).
- 830 [22] Dorigi, K. M. *et al.* Mll3 and Mll4 Facilitate Enhancer RNA Synthesis and Transcription from Promoters
831 Independently of H3K4 Monomethylation. *Molecular Cell* **66**, 568–576 (2017).
- 832 [23] Cao, K. *et al.* An Mll4/COMPASS-Lsd1 epigenetic axis governs enhancer function and pluripotency
833 transition in embryonic stem cells. *Science Advances* **4**, eaap8747 (2018).
- 834 [24] Rapino, F. *et al.* C/EBP α Induces Highly Efficient Macrophage Transdifferentiation of B Lymphoma
835 and Leukemia Cell Lines and Impairs Their Tumorigenicity. *Cell Reports* **3**, 1153–1163 (2013).
- 836 [25] DI Tommaso, P. *et al.* Nextflow enables reproducible computational workflows. *Nature Biotechnology*
837 **35**, 316–319 (2017).
- 838 [26] Frankish, A. *et al.* GENCODE reference annotation for the human and mouse genomes. *Nucleic Acids*
839 *Research* **47**, D766–D773 (2019).
- 840 [27] Ninova, M., Tóth, K. F. & Aravin, A. A. The control of gene expression and cell identity by H3K9
841 trimethylation. *Development* **146**, dev.181180 (2019).
- 842 [28] Pervouchine, D. D. *et al.* Enhanced transcriptome maps from multiple mouse tissues reveal evolution-
843 ary constraint in gene expression. *Nature Communications* **6**, 1–11 (2015).
- 844 [29] Ernst, J. & Kellis, M. ChromHMM: Automating chromatin-state discovery and characterization. *Nature*
845 *Methods* **9**, 215–216 (2012).
- 846 [30] Hoffman, M. M. *et al.* Unsupervised pattern discovery in human chromatin structure through genomic
847 segmentation. *Nature Methods* **9**, 473–476 (2012).

- 848 [31] Song, J. & Chen, K. C. Spectacle: Fast chromatin state annotation using spectral learning. *Genome*
849 *Biology* **16**, 33 (2015).
- 850 [32] Zhang, Y., An, L., Yue, F. & Hardison, R. C. Jointly characterizing epigenetic dynamics across multiple
851 human cell types. *Nucleic Acids Research* **44**, 6721–6731 (2016).
- 852 [33] Zhang, Y. & Hardison, R. C. Accurate and reproducible functional maps in 127 human cell types via
853 2D genome segmentation. *Nucleic Acids Research* **45**, 9823–9836 (2017).
- 854 [34] Waddington, C. H. The epigenotype. *Endeavour* **1**, 18–20 (1942).
- 855 [35] Berger, S. L., Kouzarides, T., Shiekhhattar, R. & Shilatifard, A. An operational definition of epigenetics.
856 *Genes and Development* **23**, 781–783 (2009).
- 857 [36] Bannister, A. J. & Kouzarides, T. Regulation of chromatin by histone modifications. *Cell Research* **21**,
858 381–395 (2011).
- 859 [37] Rivera, C. M. & Ren, B. Mapping human epigenomes. *Cell* **155**, 39–55 (2013).
- 860 [38] Sekhon, A., Singh, R. & Qi, Y. DeepDiff: DEEP-learning for predicting DIFFerential gene expression
861 from histone modifications. *Bioinformatics* **34**, i891–i900 (2018).
- 862 [39] Yin, Q., Wu, M., Liu, Q., Lv, H. & Jiang, R. DeepHistone: A deep learning approach to predicting
863 histone modifications. *BMC Genomics* **20**, 193 (2019).
- 864 [40] Henikoff, S. & Shilatifard, A. Histone modification: Cause or cog? *Trends in Genetics* **27**, 389–396
865 (2011).
- 866 [41] Morgan, M. A. & Shilatifard, A. Reevaluating the roles of histone-modifying enzymes and their associ-
867 ated chromatin modifications in transcriptional regulation. *Nature Genetics* **52**, 1271–1281 (2020).
- 868 [42] Rickels, R. *et al.* Histone H3K4 monomethylation catalyzed by Trr and mammalian COMPASS-like
869 proteins at enhancers is dispensable for development and viability. *Nature Genetics* **49**, 1647–1653
870 (2017).
- 871 [43] Krogan, N. J. *et al.* The Paf1 complex is required for histone H3 methylation by COMPASS and Dot1p:
872 Linking transcriptional elongation to histone methylation. *Molecular Cell* **11**, 721–729 (2003).
- 873 [44] Di Tullio, A., Vu Manh, T. P., Schubert, A., Månsson, R. & Graf, T. CCAAT/enhancer binding protein
874 α (C/EBP α)-induced transdifferentiation of pre-B cells into macrophages involves no overt retrodiffer-
875 entiation. *Proceedings of the National Academy of Sciences of the United States of America* **108**,
876 17016–17021 (2011).
- 877 [45] Simpson, E. H. The Interpretation of Interaction in Contingency Tables. *Journal of the Royal Statistical*
878 *Society. Series B (Methodological)* **13**, 238–241 (1951).

- 879 [46] Wiencke, J. K., Zheng, S., Morrison, Z. & Yeh, R. F. Differentially expressed genes are marked by
880 histone 3 lysine 9 trimethylation in human cancer cells. *Oncogene* **27**, 2412–2421 (2008).
- 881 [47] Burton, A. *et al.* Heterochromatin establishment during early mammalian development is regulated
882 by pericentromeric RNA and characterized by non-repressive H3K9me3. *Nature Cell Biology* **22**,
883 767–778 (2020).
- 884 [48] Rybtsova, N. *et al.* Transcription-coupled deposition of histone modifications during MHC class II gene
885 activation. *Nucleic Acids Research* **35**, 3431–3441 (2007).
- 886 [49] Tyner, C. *et al.* The UCSC Genome Browser database: 2017 update. *Nucleic Acids Research* **45**,
887 D626–D634 (2017).
- 888 [50] Bussmann, L. H. *et al.* A Robust and Highly Efficient Immune Cell Reprogramming System. *Cell Stem*
889 *Cell* **5**, 554–566 (2009).
- 890 [51] Xie, H., Ye, M., Feng, R. & Graf, T. Stepwise reprogramming of B cells into macrophages. *Cell* **117**,
891 663–676 (2004).
- 892 [52] Dobin, A. *et al.* STAR: Ultrafast universal RNA-seq aligner. *Bioinformatics* **29**, 15–21 (2013).
- 893 [53] Li, B. & Dewey, C. N. RSEM: Accurate transcript quantification from RNA-Seq data with or without a
894 reference genome. *BMC Bioinformatics* **12**, 323 (2011).
- 895 [54] Bolstad, B. M., Irizarry, R. A., Åstrand, M. & Speed, T. P. A comparison of normalization methods
896 for high density oligonucleotide array data based on variance and bias. *Bioinformatics* **19**, 185–193
897 (2003).
- 898 [55] Nueda, M. J., Tarazona, S. & Conesa, A. Next maSigPro: Updating maSigPro bioconductor package
899 for RNA-seq time series. *Bioinformatics* **30**, 2598–2602 (2014).
- 900 [56] Marco-Sola, S., Sammeth, M., Guigó, R. & Ribeca, P. The GEM mapper: Fast, accurate and versatile
901 alignment by filtration. *Nature Methods* **9**, 1185–1188 (2012).
- 902 [57] Zhang, Y. *et al.* Model-based Analysis of ChIP-Seq (MACS). *Genome Biology* **9**, R137 (2008).
- 903 [58] Cuscó, P. & Fillion, G. J. Zerone: A ChIP-seq discretizer for multiple replicates with built-in quality
904 control. *Bioinformatics* **32**, 2896–2902 (2016).
- 905 [59] Landt, S. G. *et al.* ChIP-seq guidelines and practices of the ENCODE and modENCODE consortia.
906 *Genome Research* **22**, 1813–1831 (2012).
- 907 [60] Quinlan, A. R. & Hall, I. M. BEDTools: a flexible suite of utilities for comparing genomic features.
908 *Bioinformatics* **26**, 841–842 (2010).
- 909 [61] Pohl, A. & Beato, M. bwtool: A tool for bigWig files. *Bioinformatics* **30**, 1618–1619 (2014).

- 910 [62] Visser, I. & Speekenbrink, M. depmixS4: An R package for hidden markov models. *Journal of*
911 *Statistical Software* **36**, 1–21 (2010).
- 912 [63] Gabadinho, A., Ritschard, G., Müller, N. S. & Studer, M. Analyzing and visualizing state sequences in
913 R with TraMineR. *Journal of Statistical Software* **40**, 1–37 (2011).
- 914 [64] Gu, Z., Eils, R. & Schlesner, M. Complex heatmaps reveal patterns and correlations in multidimen-
915 sional genomic data. *Bioinformatics* **32**, 2847–2849 (2016).
- 916 [65] van de Velden, M., D’Enza, A. I. & Palumbo, F. Cluster Correspondence Analysis. *Psychometrika* **82**,
917 158–185 (2017).
- 918 [66] Markos, A., D’Enza, A. I. & van de Velden, M. Beyond tandem analysis: Joint dimension reduction
919 and clustering in R. *Journal of Statistical Software* **91**, 1–24 (2019).
- 920 [67] Falcon, S. & Gentleman, R. Using GOstats to test gene lists for GO term association. *Bioinformatics*
921 **23**, 257–258 (2007).
- 922 [68] Supek, F., Bošnjak, M., Škunca, N. & Šmuc, T. REVIGO Summarizes and Visualizes Long Lists of
923 Gene Ontology Terms. *PLoS ONE* **6**, e21800 (2011).
- 924 [69] Wickham H. *ggplot2: Elegant Graphics for Data Analysis* (Springer-Verlag, New York City, New York,
925 2009).
- 926 [70] Davis, C. A. *et al.* The Encyclopedia of DNA elements (ENCODE): Data portal update. *Nucleic Acids*
927 *Research* **46**, D794–D801 (2018).
- 928 [71] Dunham, I. *et al.* An integrated encyclopedia of DNA elements in the human genome. *Nature* **489**,
929 57–74 (2012).
- 930 [72] Team R. C. *R: A Language and Environment for Statistical Computing* (R Foundation for Statistical
931 Computing, Vienna, Austria, 2017).
- 932 [73] Benjamini, Y. & Hochberg, Y. Controlling the False Discovery Rate: A Practical and Powerful Approach
933 to Multiple Testing. *Journal of the Royal Statistical Society: Series B (Methodological)* **57**, 289–300
934 (1995).

935 Figures Legends

936 **Figure 1: Global behaviour and relationship between chromatin and expression during transdiffer-**
937 **entiation** — See also Supplementary Figures 1-2, 7; Supplementary Tables 1-2. **a:** The transdifferentiation
938 of human pre-B cells into macrophages lasts a period of seven days, which we monitored at twelve time-
939 points. **b:** We have performed ChIP-seq of nine histone modifications and RNA-seq in whole-cell fraction,
940 at twelve time-points along the process of transdifferentiation. All experiments were performed in two
941 biological replicates. **c:** Trajectories of transdifferentiation derived from a Principal Component Analysis
942 performed jointly on time-series gene expression and chromatin marks' profiles. **d:** Correlations between
943 levels of gene expression and histone marks. For a given mark and for each of the twelve time-points,
944 we computed the steady-state Pearson r value between the vector of expression levels and the vector of
945 chromatin signals corresponding to the 12,248 genes. These twelve correlation values are represented by
946 single dots, the size of the dot being proportional to the hours of the corresponding time-point. The median
947 Pearson r values for each mark are: H3K27ac: 0.67; H3K9ac: 0.72; H4K20me1: 0.59; H3K36me3: 0.72;
948 H3K4me3: 0.70; H3K4me1: 0.51; H3K4me2: 0.61; H3K9me3: -0.07; H3K27me3: -0.17. In the case of
949 time-course correlations, we obtained a Pearson r value for each expressed gene, and the distributions for
950 all genes are represented by violin and box plots. Median Pearson r values across genes for each mark
951 are: H3K27ac: 0.41; H3K9ac: 0.44; H4K20me1: 0.45; H3K36me3: 0.43; H3K4me3: 0.29; H3K4me1:
952 0.10; H3K4me2: 0.10; H3K9me3: 0.13; H3K27me3: -0.03.

953 **Figure 2: Genes are characterized by a limited number of major chromatin states, which are**
954 **more stable than expression** — See also Supplementary Figure 3. **a:** A five-state multivariate HMM.
955 Each state is defined by a combination of histone marks. We report the histone marks' signals corre-
956 sponding to each state. The states are sorted by increasing level of marking averaged over the nine
957 histone modifications, with a and e states characterized by the lowest and highest average level of mark-
958 ing, respectively. **b:** Heatmap representing the hierarchical clustering of the HMM profiles built along the
959 transdifferentiation process for the 12,248 genes. **c:** Arc diagram representing the types of state transi-
960 tions observed in the HMM-sequence profiles of DE genes. The size of the arrow base is proportional to
961 the number of genes reporting a given transition. Only transitions involving ≥ 10 genes are shown. We
962 tested, for the sets of genes reporting each type of transition, the significance in gene expression fold-
963 change (FC) (Wilcoxon Rank-Sum paired test, two-sided). The color of the arrow represents the average
964 FC among genes experiencing a given transition. Transitions characterized by no significant changes in
965 expression FC (Benjamini-Hochberg FDR ≥ 0.05) are represented by gray arrows. Upper panel: transitions
966 from weaker to stronger active chromatin marking. Lower panel: transitions from stronger to weaker active
967 chromatin marking. **d:** Examples showing different HMM states along transdifferentiation. For each gene,
968 expression and chromatin tracks from one biological replicate are displayed, as well as normalized line
969 plots averaging the signal from the two replicates. Profiles of HMM states for the three genes are shown
970 at the bottom. Left panels: example of an up-regulated gene (*NUCB1*) with a constant HMM state profile
971 along transdifferentiation. Middle panels: example of an up-regulated gene (*CD163*) transitioning first from

972 absence of marking state (a) to low marking state (b), and from this to strong marking state (e). Right
973 panels: example of a down-regulated gene (*MCAM*) transitioning from active marking state (d) to bivalent
974 marking state (c).

975 **Figure 3: Uncoupling of expression and chromatin marks throughout transdifferentiation** — See
976 also Supplementary Figure 4, Supplementary Tables 3-4. **a:** Expression and chromatin profiles across the
977 12 time-points (columns) for the set of 8,030 DE genes, distinguishing between differentially marked (DM),
978 stably marked (SM) and unmarked (UM) genes (rows). The profiles consist of row-normalized z-scores,
979 computed independently for expression and chromatin marks. **b:** Expression and chromatin profiles over
980 the 12 time-points (columns) for the set of stably expressed genes that are differentially marked for a given
981 histone modification along transdifferentiation. The profiles consist of row-normalized z-scores, computed
982 independently for expression and chromatin marks. The largest numbers of significantly variable profiles
983 are observed for H3K27ac and H3K9ac. **c:** analogous representation to Figure 3b for silent genes. In this
984 case, H3K4me1 and H3K4me2 are the most variable marks throughout the process.

985 **Figure 4: Chromatin marks show a coordinated behavior along transdifferentiation** — See also
986 Supplementary Figure 5, Supplementary Table 3. **a:** Decision-tree approach to label each of the 8,030
987 DE genes based on their chromatin marking status and its relationship with the expression profile over
988 time. The approach is applied independently for each of the nine histone marks. The first branch dis-
989 tinguishes between unmarked (absence of peaks across all twelve time-points) and marked (presence of
990 peaks in at least one time-point) genes. Within the set of marked genes, it further distinguishes between
991 stably and differentially marked genes, i.e. genes characterized by absence and presence, respectively, of
992 significant (maSigPro Benjamini-Hochberg FDR < 0.05) changes in chromatin signal along the process.
993 Differentially marked genes are further classified into genes with positive, null or negative time-course
994 correlation with expression. **b:** We assessed the overlap between sets of genes corresponding to the
995 decision-tree labels across different histone marks (hypergeometric test). Hierarchical clustering of the
996 FDR values identifies three main clusters: a) genes showing expression profiles positively correlated with
997 H3K27ac, H3K9ac, H3K4me3, H3K36me3, H3K4me1, H3K4me2, H4K20me1, and negatively correlated
998 with H3K27me3; b) genes unmarked for H3K27ac, H3K9ac, H3K4me3, H3K4me1, H3K4me2, H4K20me1
999 and H3K36me3; c) genes with stable or uncorrelated profiles for H3K27ac and H3K9ac, stable profiles
1000 for H3K4me3, H3K36me3, H3K4me1, H3K4me2, H4K20me1, and unmarked for H3K27me3. The color
1001 code for the labels is analogous to Figure 4a. **c:** Similar results are obtained with Cluster Correspondence
1002 Analysis, a method that combines dimension reduction and cluster analysis for categorical data. Three-
1003 dimensional representation of the genes (analysis objects), grouped into three clusters (color-coded) based
1004 on the combinations of histone marks and labels they display.

1005 **Figure 5: Chromatin marking is associated with expression specifically at the time of gene acti-**
1006 **vation** — See also Supplementary Figure 5, Supplementary Tables 5-6. **a:** Percent stacked bar plot rep-
1007 resenting, for each of the three clusters, the proportion of unmarked, stably marked, positively correlated,
1008 uncorrelated, and negatively correlated genes identified with respect to each histone mark. **b:** Examples of

1009 genes belonging to each cluster. For each gene, expression and chromatin tracks from one biological repli-
1010 cate are displayed, as well as normalized line plots averaging the signal from the two replicates. Profiles of
1011 HMM states for the three genes are shown at the bottom. Upper panels: example of an up-regulated gene
1012 (*ALDH3B1*) showing stable and uncorrelated profiles for active marking and absence of H3K9me3 and
1013 H3K27me3 along transdifferentiation. Middle panels: example of an up-regulated gene (*DAPP1*) showing
1014 positively correlated profiles for active marking and absence of H3K9me3 and H3K27me3 along transd-
1015 ifferentiation. Lower panels: example of a down-regulated gene (*U2AF1*) showing absence of marking
1016 along transdifferentiation. **c**: Percent stacked bar plot reporting the proportion of up-regulated genes in
1017 clusters 1-3 characterized by decreasing degrees of gene expression activation (bins of 10% decrement)
1018 at time-point 0h p.i. The degree of gene expression activation is defined as the ratio between the gene's
1019 expression level at 0h and its maximum expression level along transdifferentiation.

1020 **Figure 6: Gene expression changes anticipate changes in most active marks for up-regulated**
1021 **genes** — See also Supplementary Figure 6. **a**: Alluvial plot describing, for each of the seven canonical
1022 active histone marks, the number of genes, out of 257 genes activated during transdifferentiation (i.e. up-
1023 regulated genes not expressed (< 1 TPM) at 0 hours p.i.), for which the up-regulation in a given mark's sig-
1024 nal anticipates (light green), co-occurs with (green) or follows (dark green) gene expression up-regulation.
1025 For more details see Supplementary Figure 6b. The flow lines indicate the number of genes exchanged
1026 among the three groups across increasing degrees of up-regulation. **b**: Lag (hours) between 25% up-
1027 regulation in histone marks' signal and expression level for the 257 selected up-regulated genes. Negative
1028 lags correspond to changes in chromatin marks anticipating changes in gene expression; positive lags cor-
1029 respond to changes in chromatin marks following changes in gene expression. **c**: Upper panel: Heatmaps
1030 reporting the proportion (%) of genes activated during transdifferentiation whose changes in the chromatin
1031 mark on row i anticipate changes in the chromatin mark on column j . Like in the previous analyses, we con-
1032 sidered four subsequent degrees of up-regulation (25%, 50%, 75% and 100%). e.g. the fraction reported
1033 in cell [row 1, column 2] of the first heatmap (25%), corresponds to the percentage of genes for which
1034 the 25% up-regulation in H3K4me1 signal (yellow - row 1) anticipates the 25% up-regulation in H3K4me2
1035 signal (ochre - column 2). Lower panel: analogous to upper panel for the 629 up-regulated genes already
1036 expressed (> 25 TPM) at 0h p.i. For this latter set of genes there is not a precise order of increase in
1037 chromatin marks. **d**: Mean and standard deviation of time-series expression and chromatin profiles for the
1038 257 (left panel) and 629 (right panel) up-regulated genes that are not expressed and highly expressed,
1039 respectively, at 0 hours p.i. The expression and histone marks' time-series profiles of each gene were
1040 re-scaled to a 0-100% range prior to the analysis. We highlight in black the time-points at which the mean
1041 value is $\geq 25\%$.

1042 **Figure 7: A model to explain the coupling between transcription and chromatin marking over**
1043 **time a**: According to our model, chromatin marking correlates with expression specifically during the first
1044 stage of gene activation, and the deposition of histone marks follows a specific order. Further changes
1045 in gene expression that happen later in time are mostly uncoupled from chromatin marking. **b**: Examples

1046 of up-regulated genes inactive (*CCL2*) and highly active (*FTL*) at the beginning of transdifferentiation. For
1047 each gene, expression and chromatin tracks from one biological replicate are displayed, as well as normal-
1048 ized line plots averaging the signal from the two replicates. Profiles of HMM states for the two genes are
1049 shown at the bottom. Left panels: for *CCL2*, most active histone modifications follow gene activation, with
1050 the exception of H3K4me1 and H3K4me2, which anticipate it. Right panels: for *FTL*, most active histone
1051 modifications remain stable along transdifferentiation, even though its absolute increase in expression is
1052 much higher than that of *CCL2*. **c:** Percentage (%) of unmarked, stably marked, positively correlated, un-
1053 correlated and negatively correlated profiles within cluster 3, cluster 2 (0-25%, 25-75%, 75-100% activation
1054 level at time-point 0h), and cluster 1 up-regulated genes. Positively correlated genes are further sepa-
1055 rated into genes whose histone mark's up-regulation anticipates, co-occurs with or follows gene expression
1056 up-regulation.

Figure 1

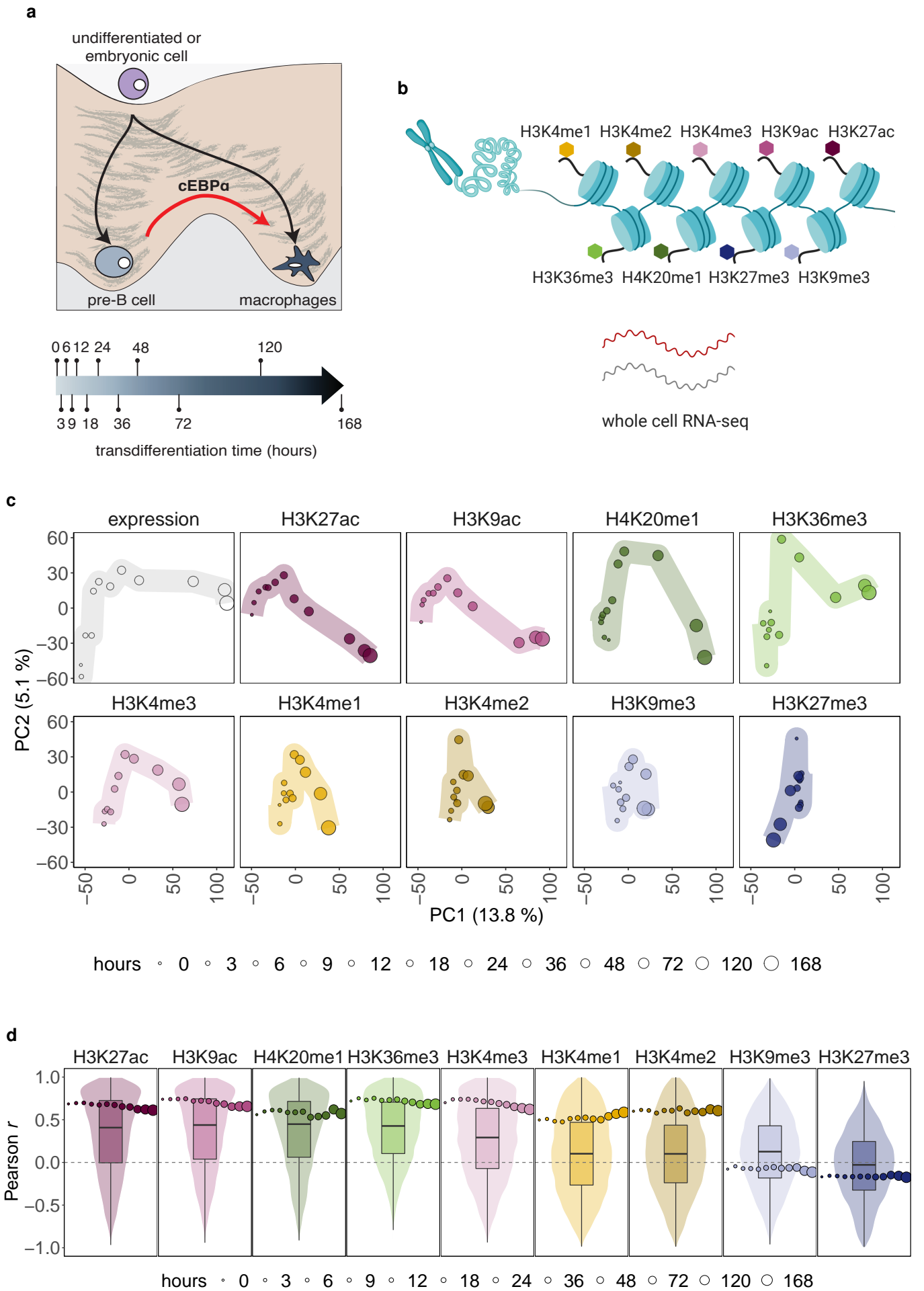


Figure 2

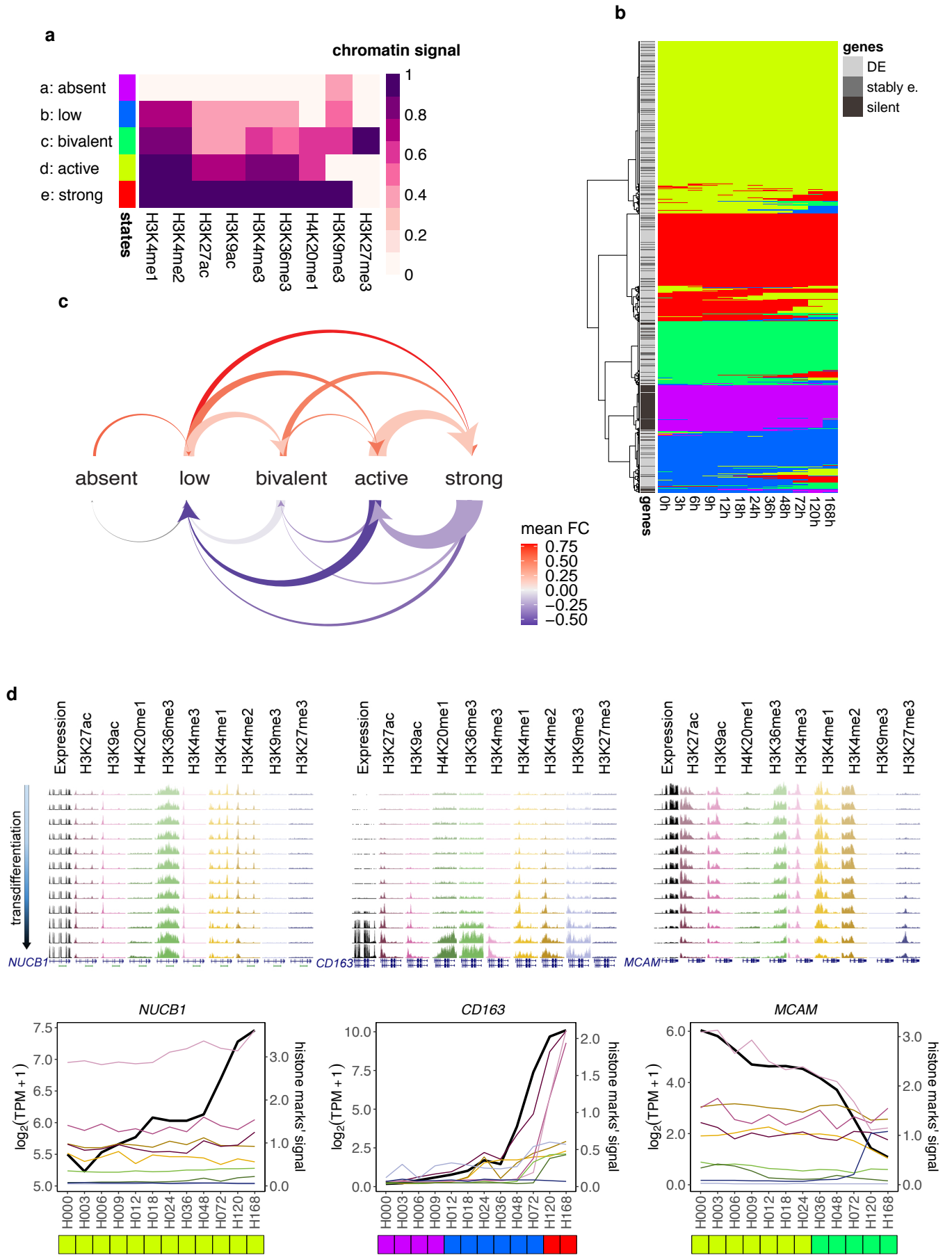


Figure 3

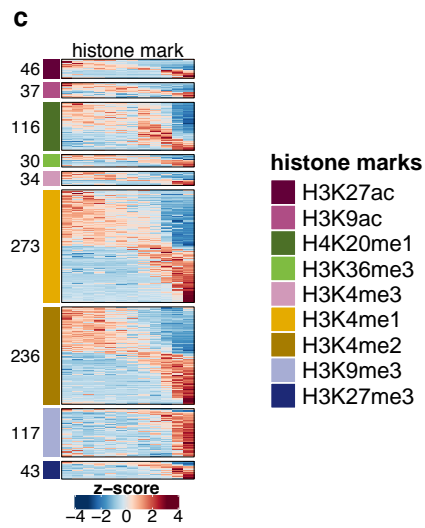
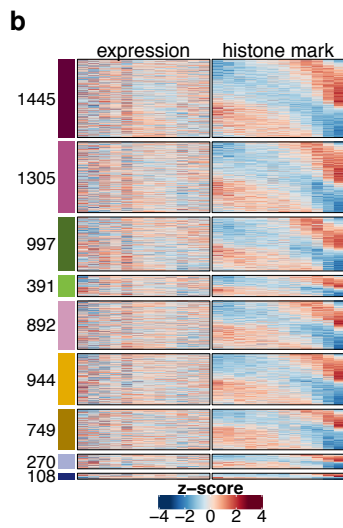
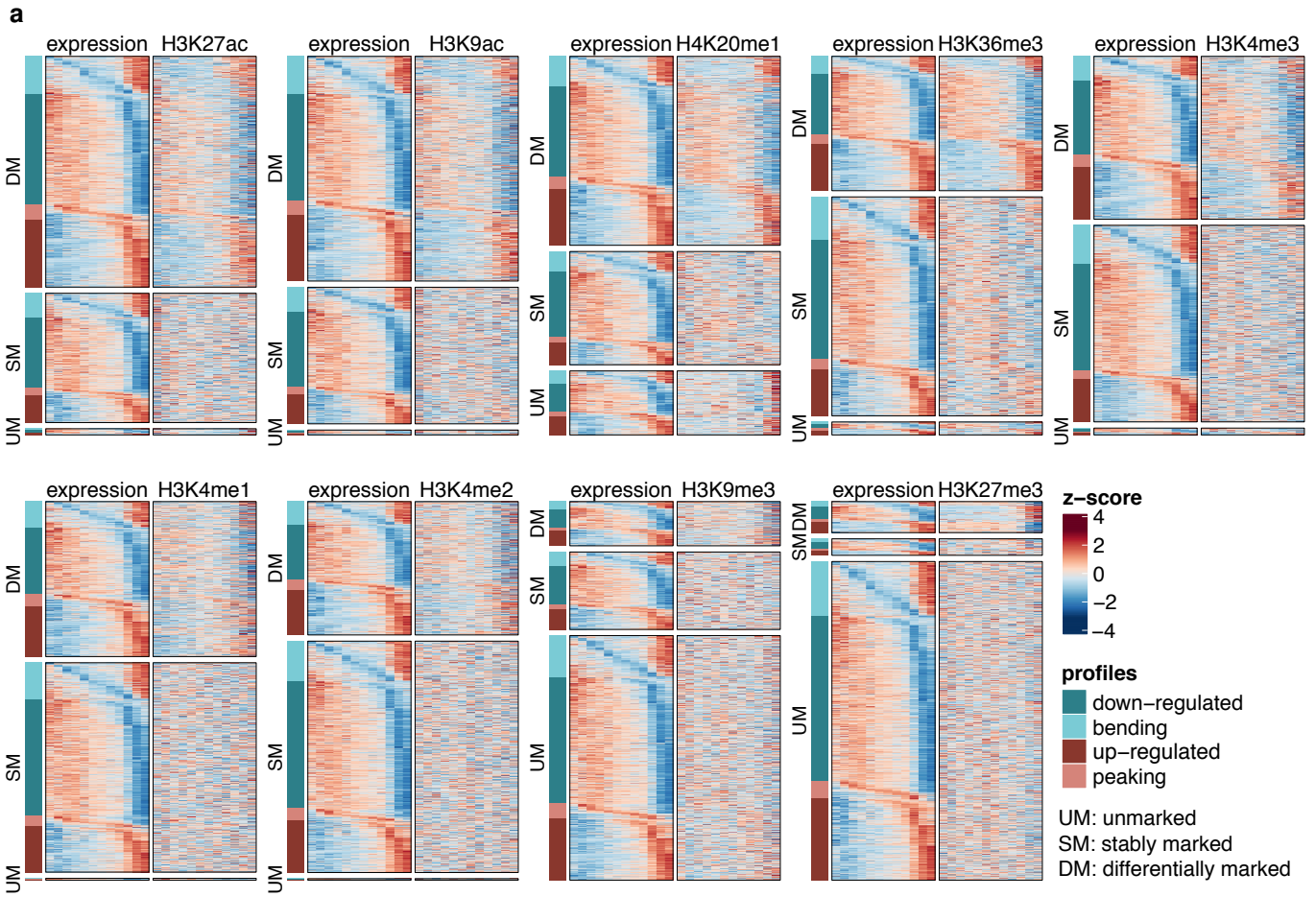
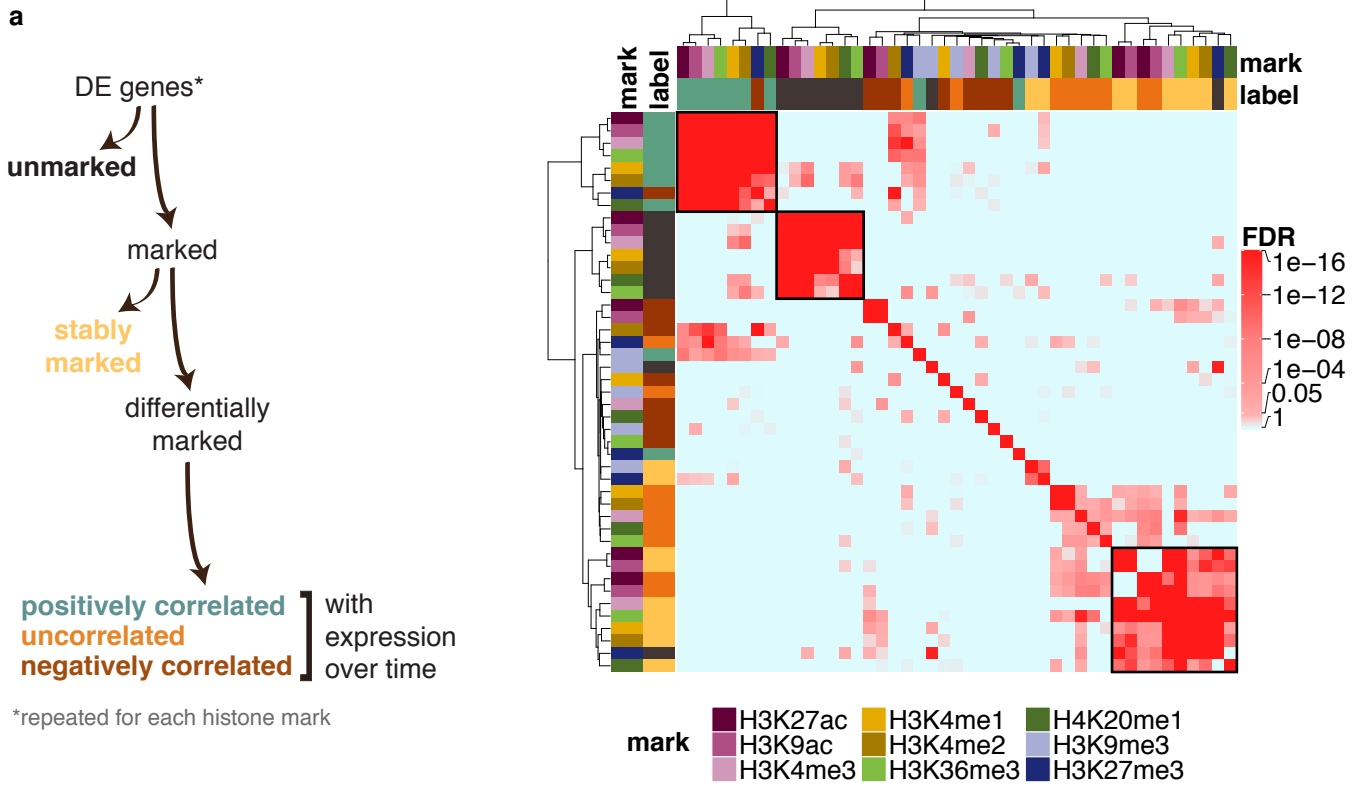


Figure 4



*repeated for each histone mark

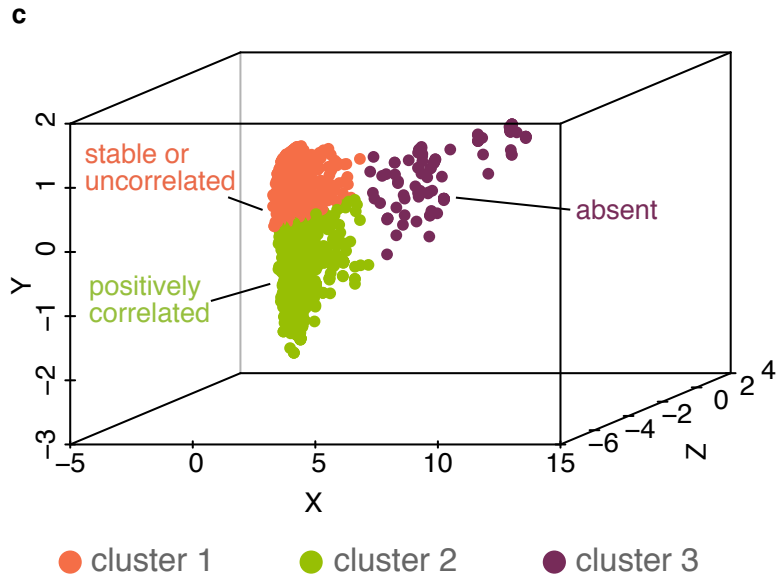
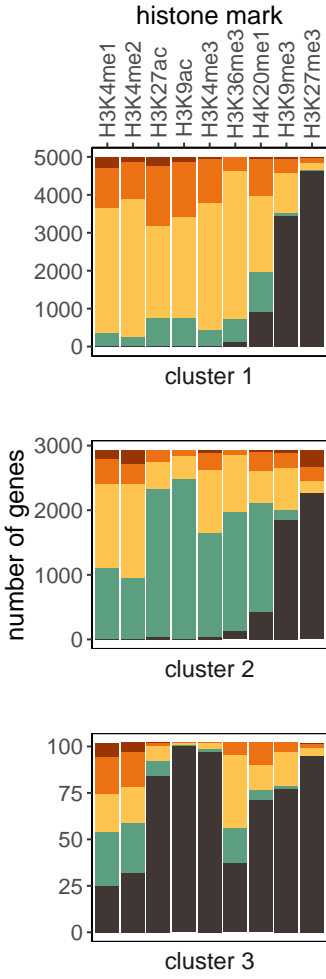
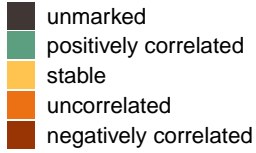


Figure 5

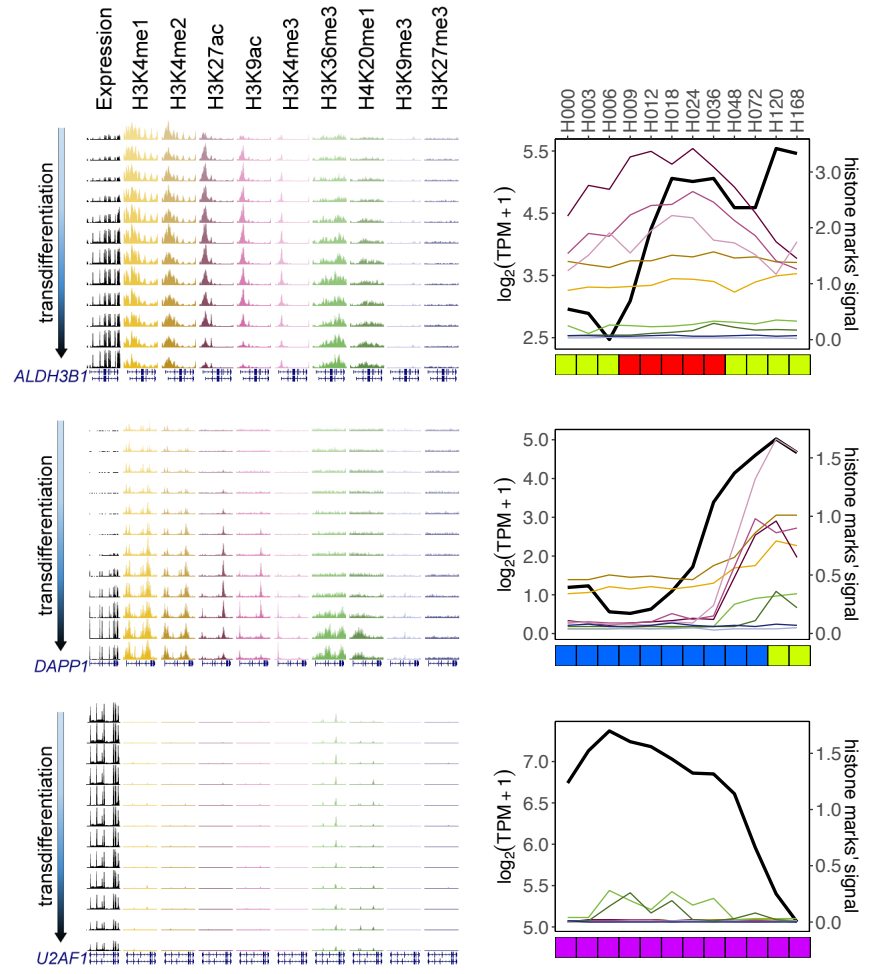
a



label



b



HMM states



c

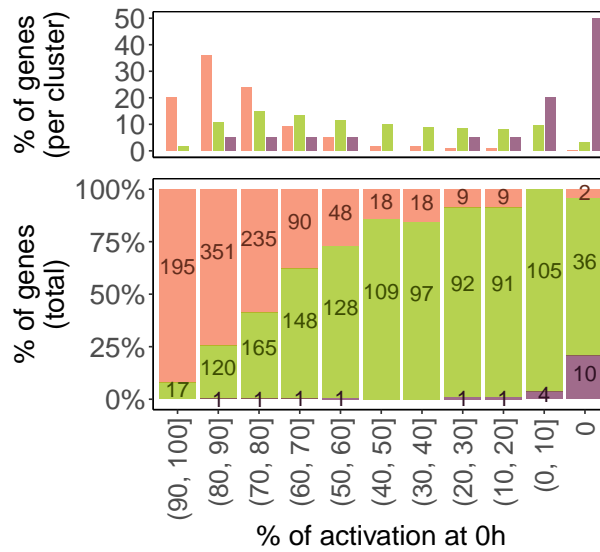


Figure 6

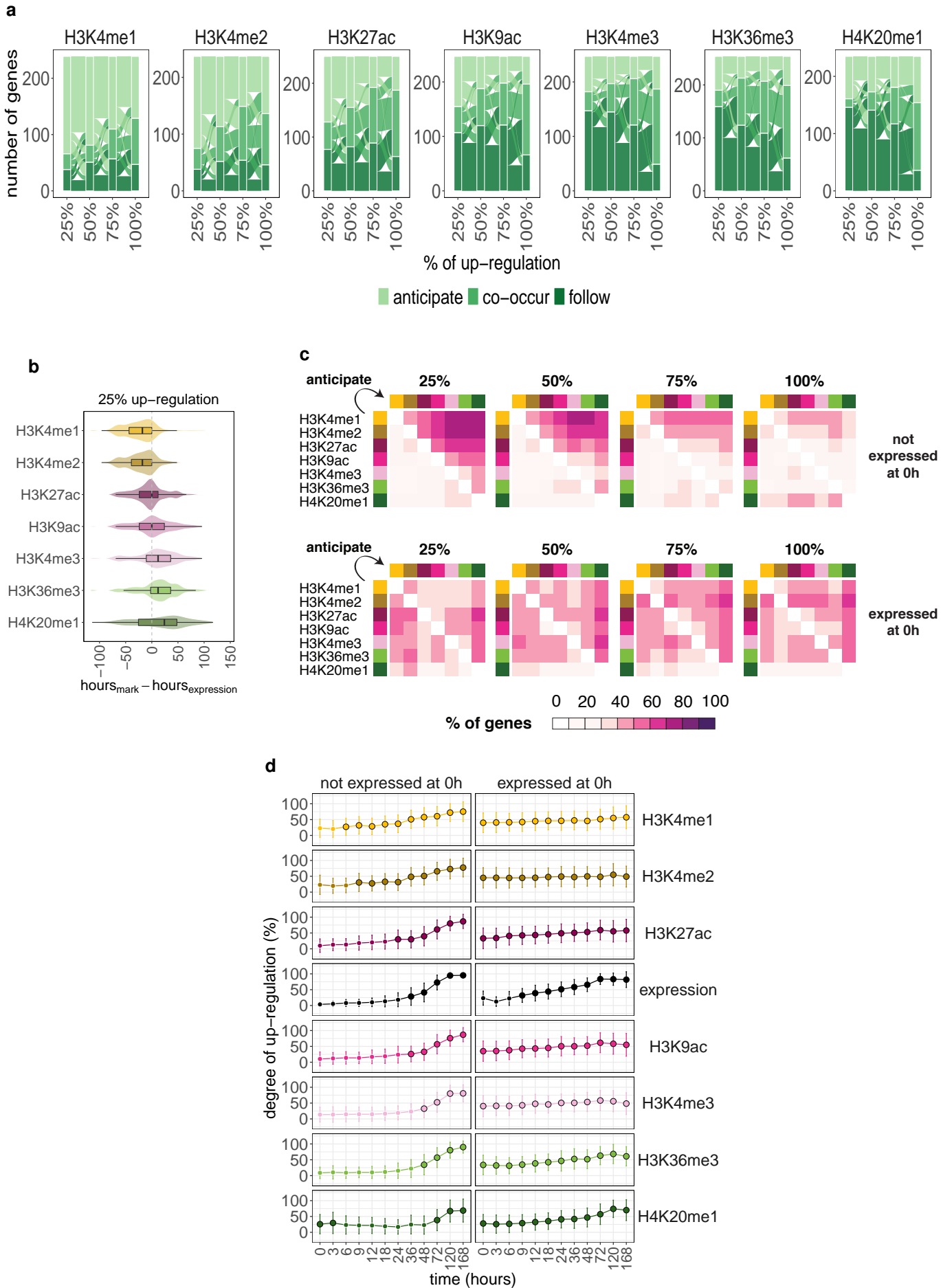


Figure 7

