

1 **Quantification of differential transcription factor activity and** 2 **multiomics-based classification into activators and** 3 **repressors: *diffTF***

4 Ivan Berest^{1,4*}, Christian Arnold^{1,*}, Armando Reyes-Palomares¹, Giovanni Palla¹, Kasper Dindler
5 Rasmussen^{2,3}, Kristian Helin^{2,3} & Judith B. Zaugg¹

6 ¹ Structural and Computational Biology Unit, European Molecular Biology Laboratory, Heidelberg

7 ² Biotech Research and Innovation Centre (BRIC), University of Copenhagen, Copenhagen

8 ³ Novo Nordisk Foundation Center for Stem Cell Biology, Copenhagen

9 ⁴ Collaboration for joint PhD degree between EMBL and Heidelberg University, Faculty of
10 Biosciences

11 * equal contribution

12 **Transcription factor (TF) activity is an important read-out of cellular signalling pathways**
13 **and thus to assess regulatory differences across conditions. However, current**
14 **technologies lack the ability to simultaneously assess activity changes for multiple TFs**
15 **and in particular to determine whether a specific TF acts globally as transcriptional**
16 **repressor or activator. To this end, we introduce a widely applicable genome-wide**
17 **method *diffTF* to assess differential TF activity and to classify TFs as activator or**
18 **repressor (available at <https://git.embl.de/grp-zaugg/diffTF>). This is done by integrating**
19 **any type of genome-wide chromatin accessibility data with RNA-Seq data and in-silico**
20 **predicted TF binding sites. We corroborated the classification of TFs into repressors and**
21 **activators by three independent analyses based on enrichments of active/repressive**
22 **chromatin states, correlation of TF activity with gene expression, and activator- and**
23 **repressor-specific chromatin footprints. To show the power of *diffTF*, we present two**
24 **case studies: First, we applied *diffTF* in to a large ATAC-Seq/RNA-Seq dataset comparing**
25 **mutated and unmutated chronic lymphocytic leukemia samples, where we identified**
26 **dozens of known (40%) and potentially novel (60%) TFs that are differentially active. We**
27 **were also able to classify almost half of them as either repressor and activator. Second,**

28 **we applied *diffTF* to a small ATAC-Seq/RNA-Seq data set comparing two cell types along**
29 **the hematopoietic differentiation trajectory (multipotent progenitors – MPP – versus**
30 **granulocyte-macrophage progenitors – GMP). Here we identified the known drivers of**
31 **differentiation and found that the majority of the differentially active TFs are**
32 **transcriptional activators. Overall, *diffTF* was able to recover the known TFs in both case**
33 **studies, additionally identified TFs that have been less well characterized in the given**
34 **condition, and provides a classification of the TFs into transcriptional activators and**
35 **repressors.**

36 **INTRODUCTION**

37 Transcription factors (TFs) are important for orchestrating coordinated and dynamic responses
38 to intra- and extracellular stimuli and regulating a multitude of biological processes. Indeed,
39 since many signaling cascades end in the activation of a particular set of TFs, observing a
40 change in overall TF activity can serve as a good read-out of signaling pathways that regulate
41 them (Kim et al., 2007). Transcriptional regulation is largely influenced by cell type specific
42 features such as cofactors, cooperative binding partners and local chromatin environment
43 (Whyte et al., 2013). Adding to this complexity, many TFs can act as transcriptional activators
44 and repressors depending on the cell type and growth condition (Han et al., 2015, 2018). Thus,
45 to correctly interpret the downstream effects of a change in abundance of a given TF, it is
46 important to understand its global mode of action within the specific context of the study.

47 TFs are typically lowly abundant proteins, which makes it difficult to detect them in proteomics
48 experiments (Kim et al., 2007; Teng et al., 2008), and even if they can be detected, their
49 abundance and activity do not necessarily correspond since TFs are highly regulated at the
50 post-translational level. On the other hand, chromatin immunoprecipitation followed by
51 sequencing (ChIP-Seq), which is the gold-standard technique for measuring genomic TF
52 binding events, provides information only for one TF at a time and does not detect global
53 changes in TF activity unless specific experimental normalisation methods are used (e.g.
54 spike-ins (Bonhoure et al., 2014)). Neither proteomics nor ChIP-Seq experiments can give any
55 insights into their mode of action. Finally, luciferase assays can measure the activity and mode
56 of action for a specific TF at a specific location and are therefore fairly low throughput (Komatsu
57 et al., 2010; Liu et al., 2011). Databases like *TRRUST* (Han et al., 2015, 2018) collect

58 annotations of regulation modes of TFs based on literature text mining and provide a
59 comprehensive resource for well-studied TF-target interactions. However, for the vast majority of
60 TFs, there is no consensus about their molecular functional mode of action. A general
61 framework for determining differential activity of TFs between conditions and classifying TFs into
62 transcriptional activators and repressors in a cell-type and condition-specific manner is currently
63 lacking.

64 Towards closing these gaps, we have developed an approach called *diffTF* to estimate global
65 changes in TF activities across conditions or cell types, and classify TFs into activators and
66 repressors based on the integration of genome-wide chromatin accessibility or histone mark
67 ChIP-Seq data with predicted TF binding sites and RNA-Seq data. We corroborated the
68 classification of TFs into repressors and activators by three independent analyses. First, we
69 showed that repressors and activators were enriched in repressive and active chromatin states,
70 respectively. Second, we confirmed that expression levels of repressors were anti-correlated
71 with their target genes while they were positively correlated with their activators. Third, we
72 obtained activator- and repressor-specific chromatin footprints based on TFs with a known
73 mode of action, and found that this agreed very well with the footprints obtained from the factors
74 as classified by *diffTF*.

75 We applied this approach to two case studies, one comparing two patient cohorts each with a
76 large number of heterogeneous samples, the second comparing two cell types along a
77 differentiation trajectory each with a small number of homogeneous samples. For the first study,
78 we obtained a large ATAC-Seq dataset of chronic lymphocytic leukemia (CLL) samples from
79 Rendeiro et al. ((Rendeiro et al., 2016)) from > 50 patients and a total of over 1 billion reads and
80 show that the quantification of differential TF activity by *diffTF* is highly robust with respect to a
81 wide range of parameter settings. We recapitulate many known TFs associated with CLL and
82 propose several novel TFs that are involved in processes related to CLL biology such as the
83 circadian clock. Furthermore, we were able to classify ~40% of these TFs (186) into activators
84 and repressors, thus reconciling some biological processes that seem driven by activators and
85 repressors at the same time. For the second case study we performed ATAC- and RNA-Seq on
86 murine multipotent progenitors (MPP) versus granulocyte-macrophage progenitors (GMP) in
87 quadruplicate. Again, with *diffTF* we were able to identify the known driver TFs of the

88 differentiation process, and we found that the majority of the highly differentially active TFs are
89 acting as transcriptional activators.
90 Finally, the approach has been successfully applied to identify TFs that are specifically
91 associated with TET2, an enzyme involved in DNA demethylation (Rasmussen et al., 2018), and
92 to identify novel driving factors in pulmonary artery hypertension (Reyes-Palomares et al., in
93 preparation).

94 RESULTS

95 Conceptual derivation of using open chromatin as read-out of differential TF activity

96 We define TF activity as the effect of a TF on the state of chromatin as measured by chromatin
97 accessibility assays (e.g. ATAC-Seq, DNase-Seq) or ChIP-Seq for active chromatin marks (e.g.
98 H3K27ac). This definition is based on our earlier work where we showed that genetic variants
99 affecting H3K27ac signal across individuals (hQTLs) can be explained by disruptions of TF
100 motifs whenever the hQTL-SNP overlaps with a TFBS (Grubert et al., 2015). Even though the
101 exact mechanisms of how changes in TF affinity translate to the chromatin level are unknown,
102 TF activity likely plays a causal role in mediating the effect of the DNA variant onto chromatin
103 marks (Liu et al., 2015). By reversing this argument, we propose to use the aggregate changes
104 in chromatin accessibility in the vicinity of putative binding sites of a TF as a read-out for its
105 change in activity (**Suppl. Fig. 1**). A similar concept has been proposed in other tools that
106 estimate TF activity based on ATAC or DHS data (Baek et al., 2017); (Schep et al., 2017).

107 Based on this concept, we developed *diffTF*, which is a computational approach to globally
108 assess differential TF activity between two conditions (basic mode; **Fig. 1a**, **Suppl. Fig. 2**) and
109 classify TFs into activators and repressors (classification mode, see below; **Fig. 1b**). It is based
110 on any data that measures active/open chromatin, putative binding sites for TFs of interest , and
111 optionally, for the classification mode only, matched RNA-Seq data. Briefly, for the basic mode, it
112 requires *in silico* TFBS that can be obtained using position weight matrices (PWMs) from a
113 database such as *HOCOMOCO* (Kulakovskiy et al., 2013) and a PWM scanning algorithm such
114 as *PWMScan* (Ambrosini et al., 2018) for all TFs, or from a database of ChIP-Seq data, such as
115 *ReMap* (Griffon et al., 2015). For each TFBS it then calculates the difference between two
116 conditions and summarizes their change in accessibility across all binding sites of a given TF. In
117 this step it also normalises for the GC content of the respective TFBS (**Suppl. Fig. 3**). The

118 significance is assessed using either an empirical or analytical procedure. The former assesses
119 the significance of the differential TF activity by comparing the real data against the distribution
120 of values obtained from permuting the sample labels. The analytical procedure, which is
121 particularly useful if the number of samples is too low for performing a reasonable number of
122 permutations, calculates the p-value explicitly based on a t-statistic and estimated variance (see
123 Methods for details and guidelines and **Suppl. Fig. 4**). In the basic mode, *diffTF* outputs
124 differential activity and p-value for each TF, which together allow the identification of a set of TFs
125 that show a significantly higher activity in one of the conditions (**Fig. 1a**).

126 **Conceptual derivation of using open chromatin and RNA to classify TFs as** 127 **transcriptional activators and repressors**

128 Surprisingly, little is known about whether a certain TF acts mostly as transcriptional repressor
129 or activator, and based on literature text mining, most TFs have been annotated multiple times
130 as both activator and repressor (Han et al., 2018) (**Suppl. Fig. 12b**) (Han et al., 2018) indicates
131 that the cell type or other external factors are important in determining a TF's main mode of
132 action. Therefore, we devised a cell-type specific and data-driven, multiomics approach to
133 classify TFs into activators and repressors within the framework of *diffTF* that can be run on top
134 of the basic mode (classification mode). Our classification framework is based on the fact that
135 increasing abundance of an activator TF results in increased transcription of its target genes
136 (and vice versa for repressors). Yet transcription is difficult to measure since a typical RNA-Seq
137 experiment measures the steady-state RNA level regulated by transcription and degradation.
138 We reasoned that measures of chromatin activity (such as accessibility) is a more direct
139 read-out for the mode of action of TFs. Based on this, we implemented an activator/repressor
140 classification scheme in *diffTF* using RNA-Seq data as an estimate for TF abundance. For each
141 TF, we calculate the correlations across individuals between its expression level and the
142 ATAC-Seq signal in its putative target peak (**Fig. 1b**). Each TF is then classified (i) as an
143 activator when it shows an overall positive correlation with the ATAC-Seq signal at its putative
144 target sites, or (ii) as a repressor for an overall negative correlation, or (iii) as undetermined if
145 the distribution of correlations is not significantly different from the peaks that did not overlap its
146 putative binding sites (see also Methods and **Suppl. Fig. 12a**). The assumptions underlying this
147 classification are tested in the context of case-study I (see below).

148

Case-study I: Quantify differential TF activity in a large ATAC-Seq dataset in CLL

149 We sought to apply *diffTF* to a large ATAC-Seq data set in a biological setting that is
150 well-studied so that we could assess the technical robustness and its power to recover relevant
151 biological signal. To do so, we identified a large ATAC-Seq data-set comparing different
152 subtypes of the extensively studied cancer chronic lymphocytic leukemia (CLL) (Rendeiro et al.,
153 2016) as an ideal dataset.

154 Chronic lymphocytic leukemia (CLL) is one of the most frequent types of cancer in the Western
155 world, particularly among the elderly. There are two major subtypes of CLL, which are defined
156 based on the mutation status of the IgHV locus (mutated: M-CLL and unmutated: U-CLL). In
157 M-CLL, B-cells go through normal affinity maturation with the aid of T-helper cells and undergo
158 multiple rounds of IgHV somatic hypermutation to produce high affinity B-cell receptors (BCR).
159 This process is essential for their development, survival and growth (Neu and Wilson, 2016). In
160 contrast, U-CLL B-cells reach their affinity maturation point in an unregulated manner, and
161 without involvement of T-helpers (Chiorazzi and Ferrarini, 2011). Overall, this leads to worse
162 clinical outcomes such as shorter survival time and higher frequency of relapse after treatment
163 (Furman et al., 2014).

164 The CLL dataset is comprised of ATAC-Seq data for a cohort of 88 CLL patients stratified by the
165 mutation status of the IgHV locus (34 U-CLL, 50 M-CLL, 4 unclassified). After data processing
166 and quality control, 25 and 27 U-CLL and M-CLL samples remained, and we applied *diffTF* in
167 the basic mode to identify the differences in TF activities between U-CLL and M-CLL (**Suppl.**
168 **Fig. 5-7**, see Supplementary Methods for more details).

169 In total we identified 67 TFs that are differentially active (FDR < 10%) between the two subtypes
170 (**Fig. 2a**; **Suppl. Table 1**). About ~41% of the differentially active TFs have previously been
171 associated with CLL and mostly (90%) agree with the reported condition (i.e., mutated or
172 unmutated; **Suppl. Table 2**), thus providing a strong biological validation of the approach. The
173 remaining 59% may represent novel candidate TFs that can advance our understanding of the
174 disease etiology in general and the biological differences of mutated and unmutated patients in
175 particular.

176 Before focusing on the biological interpretation of the specific TFs, we used this dataset to
177 assess the technical robustness of *diffTF* with respect to TF binding site predictions. First, we

178 compared the results of *diffTF* when using putative vs. ChIP-Seq validated TF binding sites
179 since predicting TF binding sites is inherently noisy and may result in many false positive sites
180 when compared to ChIP-seq experiments (Landt et al., 2012). We observed a very strong
181 correlation of the resulting TF activity differences ($r=0.81$; **Fig. 2b**, **Suppl. Fig. 8**), which
182 indicates that *diffTF* is robust with respect to false positive binding sites. Second, we assessed
183 the parameters for TF binding site predictions and found that neither the nucleotide composition
184 of the predicted binding sites for *PWMScan* (**Suppl. Fig. 9**) nor the motif database (*JASPAR* vs.
185 *HOCOMOCO*) had a strong impact on the differential TF activity ($r=0.87$, 0.62 and 0.69 ,
186 respectively, **Fig. 2c-e**). Third, we assessed whether the size of the region surrounding the TF
187 binding site from where signal was extracted (ranging from just the 7-25 bp long binding site to
188 additional 600 bp upstream and downstream) had an impact on the results. The resulting
189 differential TF activities were strongly correlated ($r>0.9$ for 50-600 bp and $r=0.76$ for the binding
190 sites only; see Supplement and **Suppl. Fig. 10**). Additional robustness tests are described in the
191 supplement.

192 We also assessed the potential of *diffTF* to detect differential TF activities in experiments with
193 little biological signal. For this, we removed high-signal regions (i.e. differentially accessible
194 peaks at 5% FDR; see Methods) and compared the resulting differential TF activities to those of
195 the full set. We found that they were very similar for both sets ($r=0.89$), thus demonstrating the
196 power of *diffTF* to capture the differential TF activities by summarising the subtle changes in
197 chromatin accessibility across many TFBS genome-wide (**Fig. 2f**).

198 Finally, we assessed the dependency of *diffTF* results on sample size and sequencing depth.
199 Intriguingly, we found the results highly congruent across a wide range of sample sizes and
200 sequencing depths, with the majority of the significant TFs from the full dataset changing in the
201 same direction in the subsampled data (**Fig. 2g**, **Suppl. Fig. 11**). Generally, the number of
202 samples appears more important than read depth, and results using the full set were consistent
203 for a coverage as low as 1-5 million processed reads per sample (see Methods). Although the
204 subsampling results are dataset-specific and difficult to generalize they provide guidelines for
205 the applicability of *diffTF* and are in line with single-cell ATAC-Seq data analysis that also show
206 robustness for low coverage at the level of genome-wide summary statistics (Mezger et al.,
207 2018).

208 In summary, these results establish the robustness of *diffTF* in quantifying differences of TF
209 activities, and demonstrate that aggregating signal across all binding sites is a powerful and

210 sensible approach to overcome limitations such as low coverage and little underlying biological
211 signal.

212

***diffTF* proposes many novel TF candidates**

213 We next focused on the biological interpretation of the differentially active TFs (FDR 10%)
214 between M-CLL and U-CLL patients (**Fig. 2a; Suppl. Table 1**). Since TF binding motifs, which
215 are the basis of *diffTF* (and any other tool that is based on predicted TF binding sites), can be
216 very similar between TFs of the same family, we decided to group TFs into TF-motif families
217 using the PWM clustering tool *Rsat* (Medina-Rivera et al., 2015); clusters available at
218 <https://bit.ly/2J9TaaK>). The resulting clusters showed overall consistent activity changes within
219 a TF family (**Fig. 3a**), with the exception of cluster 17 that can be explained by a prominent split
220 into two branches early in the clustering into NFAT and NFκB factors, which show more activity
221 in U-CLL and M-CLL, respectively (**Fig. 3c**).

222 The most active TF cluster in U-CLL is the IRF family and STAT2 (cluster 40; **Fig. 3d**), both of
223 which have been associated with disease onset and progression, and harbour several CLL
224 susceptibility loci (Arvaniti et al., 2011; Havelange et al., 2011; Slager et al., 2013). It is followed
225 by the PAX TFs (cluster 54), which affect B-cell to plasma cell differentiation (PAX5), that is
226 linked with cell survival and poor prognosis in CLL (Ghamlouch et al., 2015). Another prominent
227 set of regulators are the members of the AP-1 complex (cluster 4), which increase proliferation
228 and play an important role in driving the invasive nature of U-CLL (Mittal et al., 2013). Finally, we
229 found c-MYC, which is involved in cell proliferation and differentiation and is highly abundant in
230 U-CLL (Landau et al., 2015; Yeomans et al., 2016).

231 For M-CLL, we identified TFs that regulate and possibly reduce apoptosis, regulate cell cycle
232 and suggest normal functionality of B-cells through the classical BCR, NF-κB and Wnt signaling
233 pathways. The most active TF family in M-CLL patients is that of the POU TFs, also known as
234 Oct factors (cluster 12; **Fig 3b**), which regulate B-cell development and immunoglobulin
235 production, therefore promoting survival of the lymphoma cells (Heckman et al., 2006). This is
236 followed by the ROR factors (cluster 36), which together with Wnt5a activate NF-κB-dependent
237 survival signaling in CLL (Minami et al., 2010), and the GATA family (cluster 16), which is known
238 to prime HSCs towards the lymphoid lineage and increase self-renewal of the stem cells in CLL

239 (Kikushige et al., 2011). Other examples include the EGR family, whose motifs are enriched in
240 aberrantly hypomethylated CpG sites in CLL (Oakes et al., 2016), PPARD, which has recently
241 been linked to M-CLL through its effect on metabolic pathways in cancer cells (Li et al., 2017),
242 and members of the GLI family, which are part of the Hedgehog signaling pathway and regulate
243 apoptosis, thereby supporting survival of M-CLL cells (Kern et al., 2015).

244 Among the novel factors associated with U-CLL, we found several TFs (i.e. BMAL1, CLOCK,
245 and NR1D1) that are involved in the regulation of the circadian clock, which has recently been
246 proposed as hallmark of cancer (El-Athman and Relógio, 2018). Moreover, we found members
247 of the basic helix-loop-helix family, such as BHE40, a regulator of mitotic division (D'Annibale et
248 al., 2014), which is essential for the development of B1-a cells (Kreslavsky et al., 2017) and
249 TFAP4, TFE3 and TFEB, for which there are known cases of gene-fusions in renal cell
250 carcinoma (Kauffman et al., 2014). Another set of novel TFs more active in M-CLL are
251 associated with pathways relevant for cancer- and B-cells such as escape from apoptosis
252 (ZN784) (Kasim et al., 2017), regulation of cell cycle progression (ZBTB6) (Chevrier et al.,
253 2014), and selection of B-cells and promotion of fetal B lymphopoiesis (ARID3A) (Zhou et al.,
254 2015). The GFI1 family (cluster 35) is less active in U-CLL and their expression and activation
255 might influence and decrease rates of apoptosis in B-cells (Coscia et al., 2011).

256 In summary, these results show that *diffTF* is able to recapitulate much of the known biology of
257 the two subtypes of CLL and, in addition, identifies several more factors that are likely to be
258 implicated in the disease.

259 **Determination of the molecular function of TFs: transcriptional repressors and** 260 **activators**

261 The paragraphs above have shown that *diffTF* can identify TFs that alter their activity across
262 different types of CLL patients. However, to gain mechanistic insights into some of the
263 regulatory differences between U-CLL and M-CLL, it is crucial to know whether a TF acts as
264 activator, in which case a higher abundance would generally result in increased transcription of
265 its target genes, or repressor, in which case an increase in abundance would be accompanied
266 by decreased target gene transcription (**Fig. 1b**). To do so, we employed the classification mode
267 of *diffTF*, which integrates the ATAC-Seq data with RNA-Seq to classify TFs as activators or

268 repressors. For this, we first needed to test the global assumption that repressors and activators
269 have an opposing effect on chromatin accessibility that underlies our classification framework.
270 For activators, the expectation is that an increased TF abundance will increase the accessibility
271 at its targets sites. For a repressor, however, the relationship between abundance and
272 accessibility at its binding site is less straightforward: On the one hand the binding of the factor
273 itself will increase the accessibility locally, while on the other hand, repression is globally
274 associated with closed chromatin. To understand the effect of repressors and activators on
275 chromatin accessibility and derive general principles, we compared the accessibility footprint
276 (Tn5 insertion sites) of a well-known repressor (REST) and a well-known activator (STAT2) that
277 are active in our cell type. We observed that for the repressor REST, there is indeed an increase
278 in accessibility at its motif, which likely reflects the binding of the TF itself. Importantly, however,
279 the accessibility drops to below the genome-wide average within 10 bp from the center of the
280 motif (**Fig. 4a, bottom**). In contrast, for the activator STAT2, we observed increased chromatin
281 accessibility outside its core binding site, which only slowly drops to the genome-wide average
282 over a distance of >100 bp from the center of the motif, likely representing the effect of the TF
283 on opening the surrounding chromatin (**Fig. 4a, top**). This shows that, while there is an increase
284 in accessibility for repressors at the immediate binding site, the surrounding chromatin is highly
285 compact while it is open for the activator. This is in line with a previous observation on EGR and
286 SP4 (Baek et al., 2017) and justifies our classification approach implemented in *diffTF*.

287 Applying this reasoning to the CLL dataset, where RNA-Seq data was available for eight
288 individuals (after QC; see supplement), we were able to classify 44% of the expressed TFs as
289 either activators or repressors (**Fig 4b-f**; n=186). Among the top activators are the IRFs, which
290 are well known transcriptional activators (Yanai et al., 2012) and which showed the same
291 footprint pattern as STAT2 (**Fig. 4e**). Among the top repressors, we found PAX5, which has
292 been shown to repress the activity of BLIMP-1 (Yasuda et al., 2012) and also shows a footprint
293 similar to the repressor REST (**Fig. 4b,e**). To assess the binding properties of activators and
294 repressors globally we performed an unbiased footprinting analysis for all TFs deemed
295 significant in CLL. Importantly, we found that the aggregate signal across all repressors
296 produced a footprint similar to that of REST, while the footprint for activators looks similar to
297 STAT2, again indicating that repressors and activators have very distinct open chromatin
298 footprints (**Fig. 4f**). Clustering of the footprints of the individual TFs revealed four major classes.
299 Class I is characterized by low levels of Tn5 insertions in the motif and high levels in the

300 adjacent regions and its members are mainly classified as activators. Class II comprises of TFs
301 with very low accessibility overall (mainly repressors). Class III contains TFs with high
302 accessibility at the binding site and low accessibility in the adjacent regions, mainly classified as
303 repressors but including a few activators. Finally, Class IV comprises TFs with a footprint that
304 neither resembles an activator nor a repressor (**Suppl. Fig. 13**). This clustering indicates that
305 TFs with a Class I footprint are likely classified as activators. In contrast, Class III footprints (like
306 REST) are more likely classified as repressor, even though there might be some activator TFs
307 that with a similar footprint. Overall, it seems that TF footprints correlate well with the molecular
308 mode of action of a TF as identified by *diffTF*.

309 When investigating TF families as defined above with the RSAT clusters (**Fig. 3**), we found that
310 TFs from the same PWM cluster are often classified both as activators and repressors (**Suppl.**
311 **Fig. 14**), supporting the hypothesis that the molecular function of a TF is highly variable. The
312 exceptions are cluster 40, containing mainly members of the IRF family, and cluster 17 that
313 contains both NFAT and NFkB TFs, which are mostly classified as activators. The circadian
314 regulators provide an example of why it is important to know the mode of action of a particular
315 TF: When analysing the differential TF activities, it appears as if BMAL1 is more active in M-CLL
316 while the other two TFs (CLOCK and NR1D1) are more active in U-CLL (**Fig 4d**). However,
317 since BMAL1 is an activator, while CLOCK and NR1D1 are repressors, all three circadian
318 factors are consistently more active in M-CLL, albeit with a contrary effect on their target genes.

319 To assess the validity of our repressor/activator classification, we chose three independent
320 approaches. First, we used *chromHMM* chromatin states for primary B-cells from the
321 Epigenomic Roadmap (Roadmap Epigenomics Consortium et al., 2015) to assess whether
322 activators and repressors are preferentially located in active and repressive states, respectively.
323 Indeed, we observed that the fraction of TFBS in active chromatin states was significantly larger
324 for activators than for repressors, and vice versa for heterochromatin and repressive states (**Fig.**
325 **5a**, see also **Suppl. Fig. 15**), thus corroborating our classification of their molecular mode of
326 action. Second, we assessed whether the direction in gene expression changes of TFs between
327 U-CLL and M-CLL was in agreement with their differential activity and molecular mode of action.
328 Again, our observations were in line with our expectations: activators showed a positive
329 correlation of activity and expression change ($r=0.19$, $P=0.05$) while repressors showed a
330 negative relationship ($r=-0.32$, $P=0.0033$; **Fig. 5b**). Third, we checked whether the expression of

331 target genes of a given TF changes in the same direction as its activity calculated by *diffTF*,
332 regardless of the TF's classification as activator or repressor, and again found that this was the
333 case for both activators and repressors (**Fig. 5c**, see Methods). In summary, these observations
334 provide three independent lines of evidence that our approach implemented in *diffTF* is able to
335 classify TFs globally by their mode of action. The fact that the correlations are in the expected
336 direction but not perfect are likely reflecting that TFs are also regulated on the
337 post-transcriptional level and thus show the limitation of using gene expression as a proxy for
338 the abundance of the active form of TFs.

339

Case study II: Applying *diffTF* to small scale multiomics dataset

340 To assess the applicability of *diffTF* to small datasets, we decided to apply it to the well-studied
341 mouse hematopoietic system. We generated ATAC-Seq and RNA-Seq profiles of multipotent
342 progenitor cells (MPP; Lin⁻Kit⁺Sca1⁺; CD150⁻CD48⁺), an early hematopoietic progenitor
343 population capable of supporting multilineage blood production (Sun et al., 2014), as well as the
344 more differentiated and myeloid-restricted granulocyte-monocyte progenitors (GMP;
345 Lin⁻Kit⁺Sca1⁺; CD16/32⁺). The profiles obtained were processed using an in-house ATAC-Seq
346 pipeline and *diffTF* (using the analytical procedure due to the small number of samples) to
347 identify TFs that are differentially active between MPP and GMP cells (see Online Methods).
348 Due to the large number of significant TFs, reflecting the high diversity between the two cell
349 types, we used RNA-Seq data to filter out non-expressed TFs. The differential signal is
350 dominated by an increase activity of the members of the well-known class of master regulators
351 of myelopoiesis, the CEBP family (C/ebp α , β , δ , ϵ , γ) in GMPs (**Fig. 6a**, **Suppl. Fig. 16**). In
352 addition, we observed higher activity of the MYC/MYB factors, which are known to be
353 exclusively active in the GMPs (Baker et al., 2014) and in NFIL3, which is involved in the
354 generation of natural killer cells (Gascoyne et al., 2009). Conversely, MPPs show a higher
355 activity for IRF/STAT, ZEB1 and ITF2 (part of the Wnt signaling) as well as TFs from the
356 Homeodomain (HXB7,HXA10) and Forkhead (FOXO3) families, all of which are associated with
357 self-renewal of hematopoietic stem cells (Sands et al., 2013).

358 The small number of samples made the correlation-based classification of the TFs into
359 activators and repressors unreliable. Therefore, we devised a second - less quantitative -
360 approach for activator/repressor classification that is based on the TF footprint and differential

361 RNA-Seq expression. In short, we determined whether TF activity and expression level co-vary
362 in the same direction and combined this with visual inspection of the pattern of its footprint
363 (**Suppl. Fig. 17**). This allowed us to identify the set of activators that showed a clear activating
364 footprint as observed for the Class I TFs in the CLL data (**Fig. 6b**). Similar to the Class II and III
365 for CLL, the pattern was less clear for the repressor footprint clusters, which contain both
366 potential repressors and activators. Interestingly, the most differentially active TFs between MPP
367 and GMP are mainly classified as activators (CEBPs, NFIL3, IRFs) or have mixed evidence (i.e.
368 activator footprint, but inconsistent directionality of expression and activity such as DBP and
369 HLF). The most differentially active repressor we identified is JUN, whose difference in activity is
370 far below the activators, indicating that the differentiation process from MPP to GMP is mainly
371 driven by transcriptional activators.

372 Overall, these results show that *diffTF* is able to identify the known TFs that drive the
373 differentiation from MPP to GMP, thus demonstrating its power to detect signals also for a small
374 number of samples. Furthermore, we show how a qualitative classification scheme of TFs into
375 activators and repressors that is primarily based on TF footprints can be useful in comparisons
376 where the small number of samples does not allow a correlation-based classification.

377

Comparison with similar tools

378 We compared the few tools with a similar focus (Baek et al., 2017; Heinz et al., 2010; Schep et
379 al., 2017) with *diffTF* (**Suppl. Table 5**) both qualitatively (*chromVAR*, *BagFoot*, *HOMER*) and
380 quantitatively (*HOMER* and *chromVAR*). Overall, *diffTF* provides a more flexible and tailored
381 analysis framework due to the extensive choice of parameters, diagnostic plots, TFBS-specific
382 results, visualizations, and pipeline adjustability. As mentioned above, it is unique in its ability to
383 directly integrate RNA-Seq with ATAC-Seq data to classify TFs into activators and repressors.
384 Due to its flexibility, *diffTF* is computationally expensive, and we provide detailed instructions on
385 memory and time requirements in the documentation.

386 We first compared *diffTF* with a more traditional TF motif analysis such as *HOMER* (Heinz et al.,
387 2010), which looks at motif enrichment in a set of differentially accessible peaks. Strikingly, no
388 enriched motifs were found in M-CLL with *HOMER*, while the few discovered in U-CLL
389 correlated significantly with differential TF activity as computed by *diffTF* (**Suppl. Fig. 18**, see

390 Methods). This analysis highlights the power of *diffTF* to capture more signal than standard
391 motif enrichment approaches.

392 To compare *diffTF* with an approach that is also based on TF activities we chose *chromVAR* and
393 *BaGFoot* (Baek et al., 2017). We were unfortunately unable to run and adopt the BagFoot
394 workflow for our CLL data due to missing example files, an incomplete documentation and
395 unresponsiveness from the authors. For *chromVAR*, the results correlate very well overall, with
396 correlation coefficients between 0.75 and 0.93 (Pearson) and 0.69-0.88 (Spearman), depending
397 on the set of TFs (i.e., all TFs or only those deemed significant by *diffTF*, therefore
398 predominantly removing TFs with low signal) and whether *chromVAR* deviations or deviation
399 scores are compared against (**Suppl. Fig. 19a-b**, see Methods). Differences likely arise due to
400 distinct methodological divergences such as comparing fold-changes for peaks (*chromVAR*) vs.
401 binding sites (*diffTF*) or whether to compare the TF-specific effect against the mean effect
402 across all TFs (*diffTF*) or not (*chromVAR*; see also Methods and **Suppl. Fig. 19c-f**). However,
403 *diffTF* goes one step beyond the currently available methods by classifying TFs based on their
404 mode of regulation - activator or repressor, thus providing important additional insights into their
405 molecular function.

406 DISCUSSION

407 We presented a genome-wide method for quantifying differences in differential TF activity for a
408 large set of TFs simultaneously, and for classifying them into their molecular mode of action as
409 transcriptional activators or repressors. The method is available for download at
410 <https://git.embl.de/grp-zaugg/diffTF> along with a comprehensive documentation and example
411 data.

412 We have shown in two case studies that *diffTF* is able to recover a change in activity for the TFs
413 expected to drive the biological processes, thus demonstrating the biological validity of the
414 method. In addition, we have extensively tested and demonstrated the technical robustness of
415 *diffTF*. In particular, we have shown that *diffTF* is able to overcome the inherent noisiness of TF
416 binding site predictions by aggregating data across all putative binding sites.

417 Calculating differential TF activity based on aggregating signal across the genome has been
418 proposed before based on the expression of putative target genes of a certain TF (Boorsma et
419 al., 2008; Bussemaker et al., 2001). Using the effect on chromatin instead of expression has
420 several advantages: first chromatin is a much simpler trait since gene expression is the sum of
421 transcription and degradation, thus increasing the power to detect differences. Second, there
422 are much more peaks than genes, thus allowing for better statistics and signal to noise ratio.
423 Finally, the effect on chromatin is much more local than on gene expression – in particular in
424 mammalian genomes, where genes can be regulated by distal enhancers. We have compared
425 differential TF activity calculated based on the average expression change of the target genes to
426 the output of *diffTF* and found that while the direction of activity between both methods is highly
427 correlated, the signal is much lower when using the expression of target genes instead of
428 chromatin accessibility at putative binding sites.

429 To demonstrate the power of *diffTF* for large but heterogeneous datasets, we have applied it to
430 identify and characterise differences between M-CLL and U-CLL from a publicly available
431 dataset of ATAC-Seq (1bn reads, 52 patients) and RNA-Seq. It is noteworthy that a TF motif
432 enrichment analysis on the significantly differentially accessible peaks did not reveal any factor
433 to be significantly enriched in M-CLL, indicating that in this case (as probably in many
434 patient-control studies) the key TFs are not necessarily switching their target enhancers and
435 promoters on and off, but rather mis-regulating many regions to a lesser extent. The advantage
436 of *diffTF* is that it can detect a slight shift in activity of the TF even if the signal differences at
437 each binding site are very low and rarely significant. It does so by averaging across all of a TFs
438 putative binding sites and is therefore more powerful than conventional enrichment analyses.

439 We have devised an approach within the *diffTF* framework to classify TFs into activators and
440 repressors based on the correlation of their expression level (RNA-Seq) and the activity of their
441 putative binding sites (ATAC-Seq). This information is highly relevant when interpreting the
442 action of TFs since it is important to know whether an upregulation of a TF would have a
443 positive or negative effect on chromatin (and therefore transcriptional) activity. Notably, this
444 classification could work even for datasets for which insufficient RNA-Seq data are available –
445 as we have shown for the MPP-GMP case – by jointly investigating TF footprints, differential
446 expression of the TF and differential TF activity.

447 It is important to note that TFs can often act as activators and repressors at different genomic
448 loci e.g. depending on their cofactors, whereas here we predict their main mode of action based
449 on the mean effect across all their predicted binding sites, and thereby lose any information
450 about bifunctionality. Furthermore, since the classification is based on correlations, it is heavily
451 dependent on variation in the RNA-Seq signal across individuals, thus biasing the TFs that can
452 be classified towards those that are variable across individuals. As a consequence, TFs whose
453 post-transcriptional regulation is not reflected in their transcript abundance will not get classified
454 correctly. Another potential misclassification may happen because of the similarity of PWMs
455 within a cluster, which makes it difficult to distinguish the exact effect of one TF while its
456 expression level is uniquely defined. As an example we cite PRDM1, which is part of IRF family
457 (cluster 40) and classified as very strong repressor, its footprint however looks more similar to a
458 typical activator (data not shown), suggesting that it is not PRDM1 driving the ATAC-Seq signal
459 in this case, but the IRFs. Thus, for distinguishing the functional roles of TFs from the same
460 cluster/family further biochemical experiments will be needed. Despite these potential pitfalls,
461 *diffTF* provides unique insights into the molecular mechanism of TFs on a global level.

462 Since many ATAC-seq experiments have a rather low number of samples, we also assessed the
463 power of *diffTF* to uncover biology in small (but more controlled) experiments. In particular, we
464 have performed a *diffTF* analysis to compare murine MPP and GMP (4 replicates each). Again
465 we identified the major TFs driving the differentiation, and were able to qualitatively classify TFs
466 into activators and repressors - in a correlation-independent approach. This classification
467 revealed that the bulk of the change in chromatin accessibility during the differentiation from
468 MPP to GMP is driven by activators. This case-study demonstrates the applicability of *diffTF* to
469 small-scale data.

470 While similar methods have been proposed for analysing ATAC-Seq data (Baek et al., 2017;
471 Schep et al., 2017), our method has several advantages when dealing with bulk ATAC-Seq data
472 and can also be used for histone mark ChIP-Seq data: (i) Unlike other methods that calculate
473 the background theoretically based on the genome-wide read depth, *diffTF* is insensitive to
474 sequence and locus dependent biases since we calculate a fold-change between conditions for
475 each region, thus normalizing for local read depth biases. This is particularly advantageous for
476 detecting small changes such as between two heterogeneous cohorts in patient-control studies.
477 (ii) *diffTF* allows integration with matching RNA-Seq data to classify TFs into activators and

478 repressors in a fully data-driven approach within the same analysis framework. Such
479 classifications are a significant help when interpreting the effects of up/down regulation of a
480 particular factors. (iii) *diffTF* provides the fold-change value for each TFBS which allows for easy
481 retrospective follow-up analysis, e.g. identifying the most differential regions regulated by a
482 specific set of TFs. (iv) Finally, our method might allow to analyse time-course data in an
483 additive manner by calculating the overall change of slope for each TF (see Methods).

484 Overall, with *diffTF* we present a multiomics data integration strategy of ATAC-Seq and
485 RNA-Seq data that calculates differential TF activity across conditions and classifies TFs based
486 on their molecular mode of action into activators and repressors. With this, *diffTF* can aid in
487 formulating testable hypotheses and ultimately improve the understanding of regulatory
488 mechanisms that are driving the differences in cell state on a systems-wide scale.

489 **METHODS**

490 Methods, including statements of data availability and any associated accession codes and
491 references, are available in the online version of the paper.

492 **ACKNOWLEDGMENTS**

493 We thank Bernd Klaus for help with the statistical part and EMBL for funding. A.R.P is recipient
494 of a postdoctoral fellowship granted by Fundación Ramón Areces. G.P. was supported by the
495 Otto Bayer Scholarship from Bayer Foundations.

496 **AUTHOR CONTRIBUTIONS**

497 J.Z. and I.B. conceived the study, C.A. and I.B. developed the computational framework and
498 performed the analyses, A.R.P. and G.P. contributed to the development and analysis, K.D.R.
499 performed the experiments, K.H. supervised the experiments, J.Z. supervised the study and
500 C.A., I.B., and J.Z. wrote the manuscript.

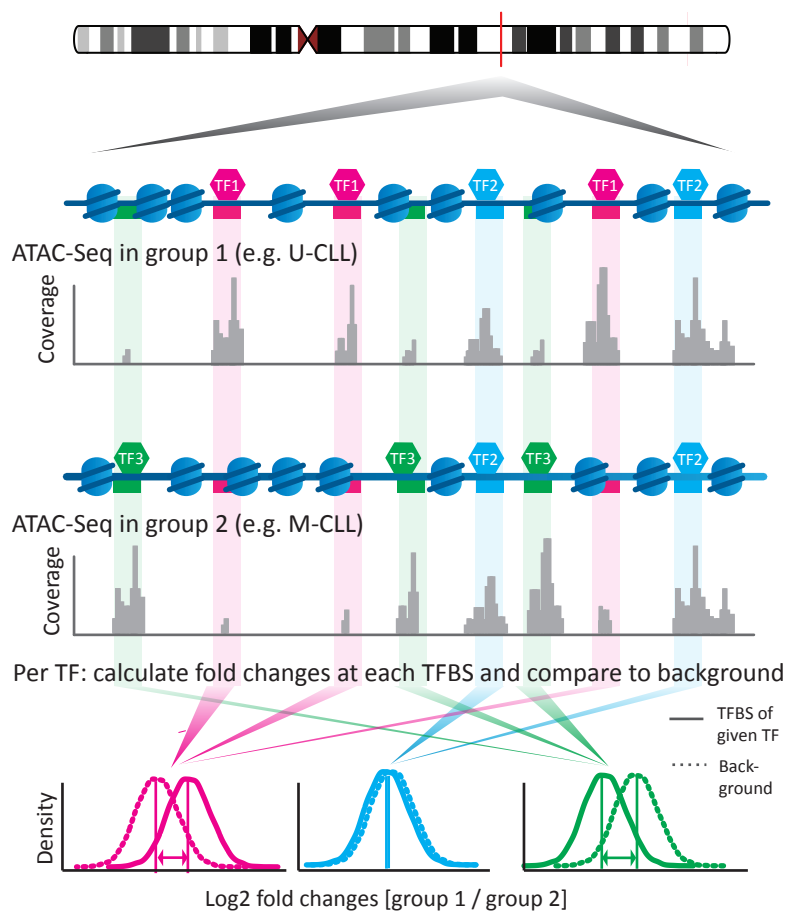
501 **COMPETING FINANCIAL INTERESTS**

502 The authors declare no competing financial interests.

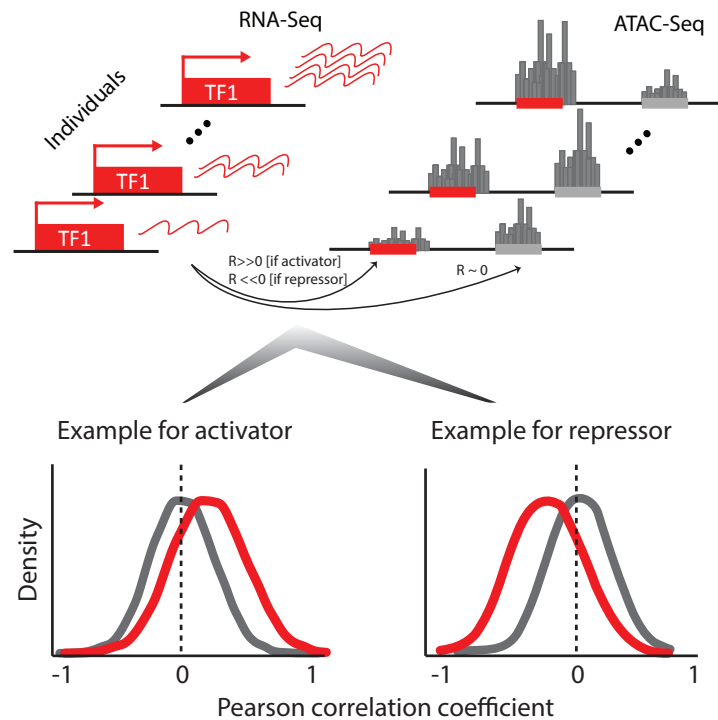
503 **Figures and figure captions**

504 **Figure 1.** *Schematic representation of the diffTF workflow. (a)* A simplified workflow illustrates
505 the principle upon which *diffTF* is based: it calculates a fold-change between two conditions for
506 each binding site of a given TF and compares this distribution to a background set of
507 fold-changes obtained from GC-content matched loci that do not contain the putative TFBS. The
508 difference in distribution is assessed in significance and effect size and visualized in a volcano
509 plot where the y-axis indicates statistical significance and the x-axis shows the effect size. (For a
510 detailed workflow see **Suppl. Fig. 1** and **2**). **(b)** *Schematic representation of the classification*
511 *approach:* correlation of TF expression level with the accessibility of its target sites. If the
512 distribution of correlations between a TF's RNA-level and the chromatin accessibility at its target
513 sites is more positive than the background distribution (accessibility at non-target sites), the TF
514 is classified as an activator in the particular biological environment; if negative, it is classified as
515 a repressor. Correlations close to 0 are classified as undetermined.

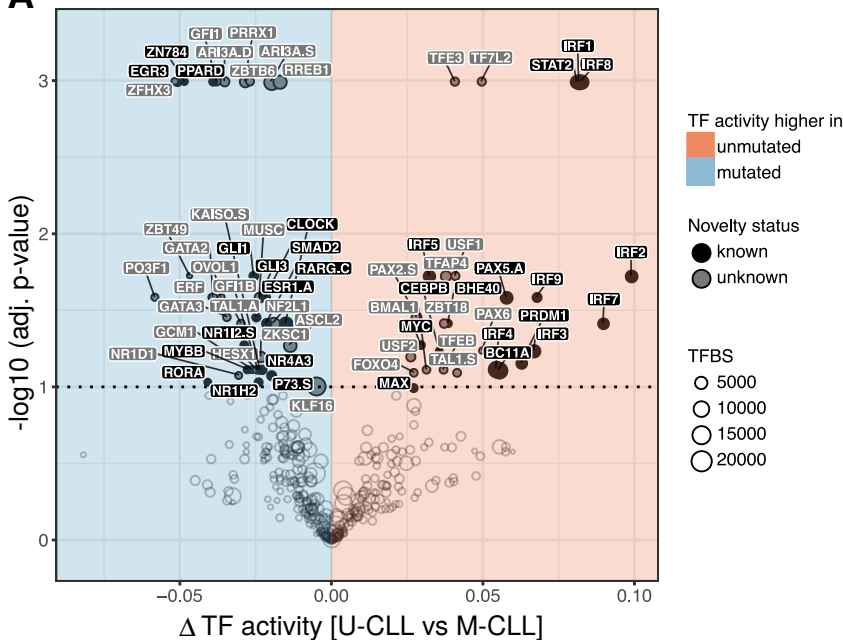
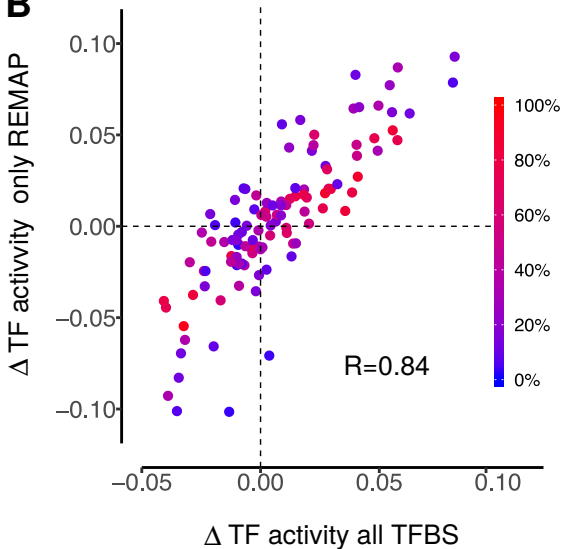
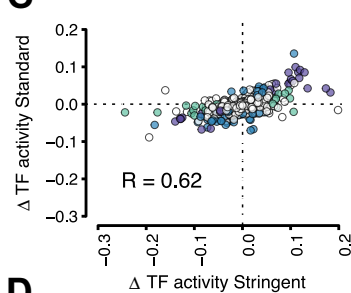
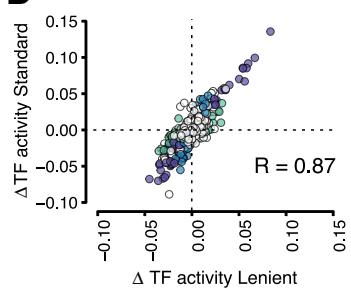
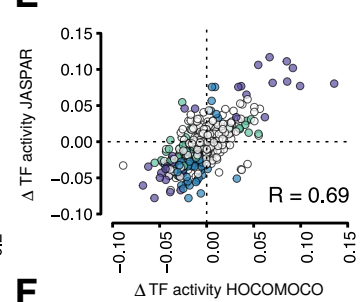
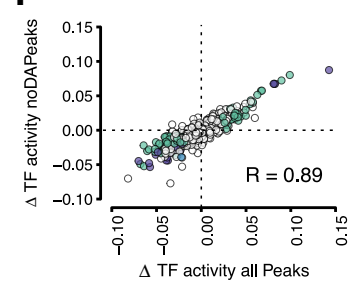
A Schematic representation of *diffTF* - basic mode



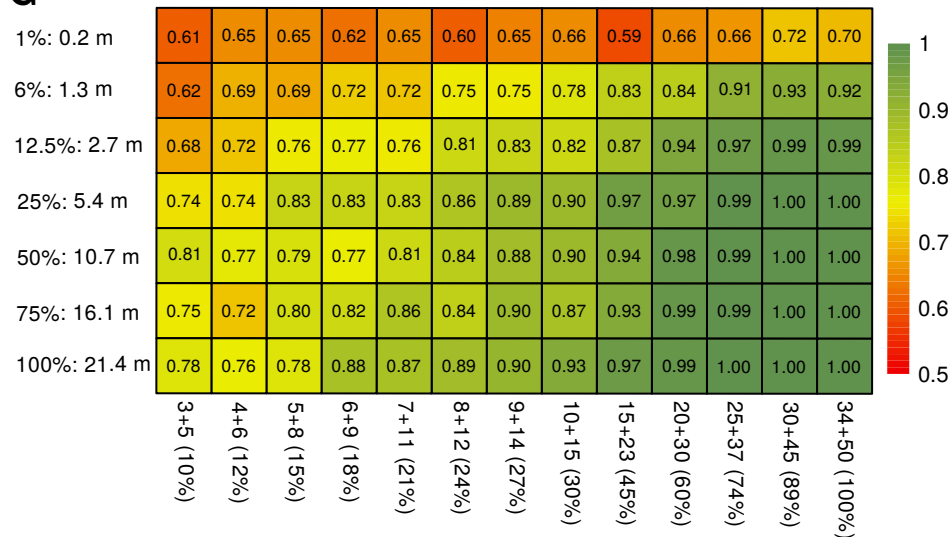
B Schematic representation of *diffTF* - classification mode



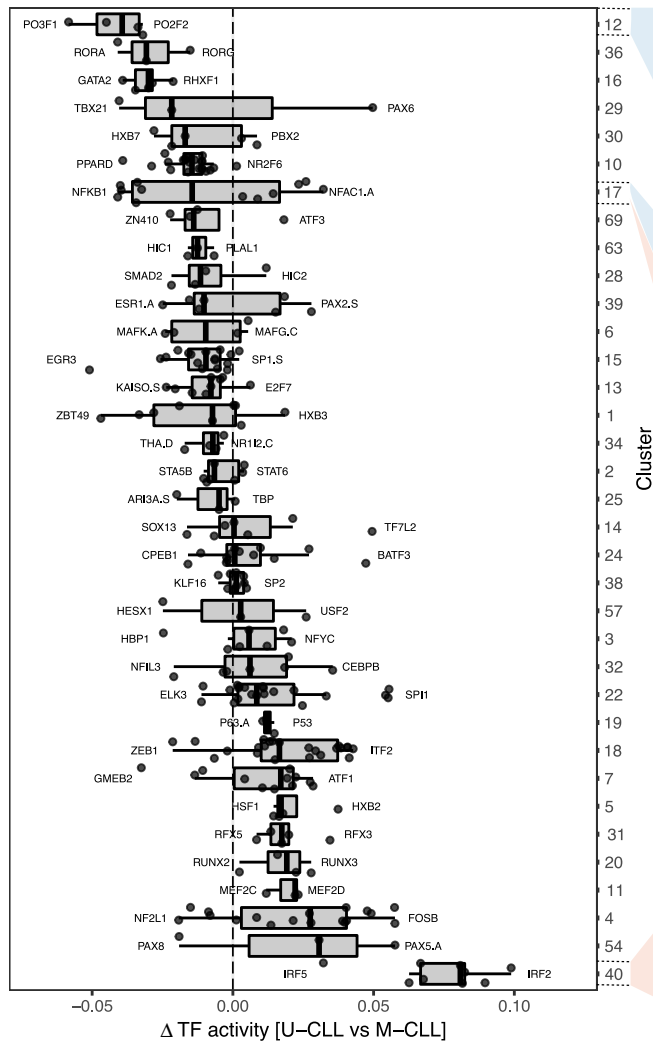
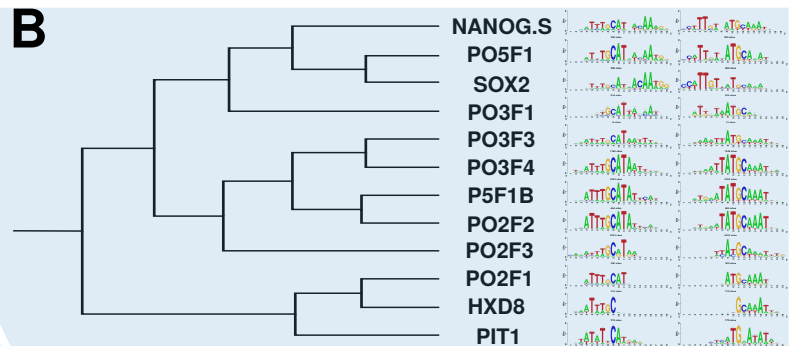
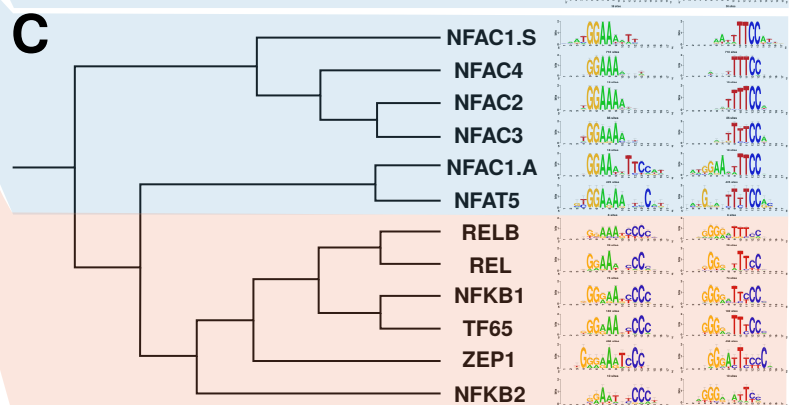
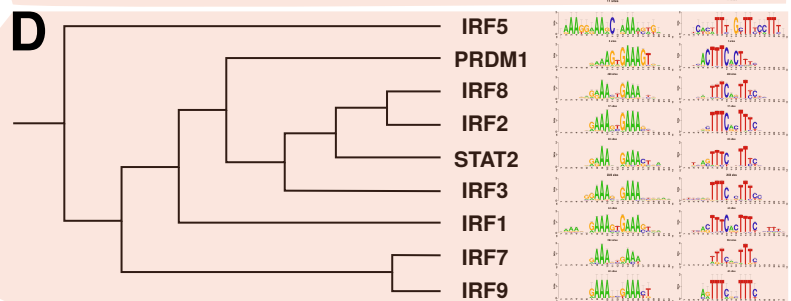
516 **Figure 2.** *diffTF* results for the CLL dataset, experimental validation and technical robustness
517 of the method. **(a)** Volcano plot of differential TF activity between U-CLL (n=27) and M-CLL
518 (n=25) patients. The y-axis denotes statistical significance (-log₁₀ transformed). TFs that pass
519 the significance threshold (5% FDR; dotted line) are labeled and colored according to their
520 novelty status (see text and **Suppl. Table. 2**). “#TFBS” denotes the number of predicted TF
521 binding sites in the peak regions for this analysis. (b)-(f): Technical robustness of *diffTF*.
522 Scatterplots of the differential TF activity from all TFs for two different *diffTF* analyses are
523 shown. Each point represents one TF. For (c-f), colors represent significance at 5% FDR (white
524 – not significant in either analysis; light green and light blue – significant for the analysis on the
525 x-axis or y-axis, respectively; purple – significant for both analyses). **(b)** Comparison of all
526 predicted TFBS and TFBS experimentally validated by ChIP-Seq data from ReMap. See also
527 **Suppl. Fig. 8.** **(c-d)** Comparison for different p-value thresholds in *PWMScan* to predict TFBS:
528 (c) standard vs. stringent (i.e., 1e-5 vs. 1e-6) and (d) standard vs. lenient (i.e., 1e-5 vs. 5e-5) for
529 a total of 628 TFs for which binding sites were retrieved for both scanning modes. **(e)**
530 Comparison of *diffTF* results based on *HOCOMOCO* v10 vs. *JASPAR* 2018 as input for the 412
531 TFs for which a motif was available in both databases. **(f)** Comparison of the full consensus
532 peak set (*allPeaks*) and only the non-differentially accessible peaks (*noDApeaks*; n=640 TFs
533 from *HOCOMOCO*). **(g)** Robustness analysis based on sequencing depth and sample size.
534 Each cell in the heatmap shows the fraction of TFs that showed the same direction of change as
535 in the full dataset for varying degrees of down-sampling sequencing depth and number of
536 samples, (5% FDR), averaged over 50 independent repetitions to minimize sampling noise.
537 Only TFs that were deemed significant in the full dataset are considered (see also **Suppl. Fig.**
538 **11**). Sequencing depth is shown as a fraction of the original data and median number of reads
539 across samples, while the number of samples is given as unmutated + mutated.

A**B****C****D****E****F**

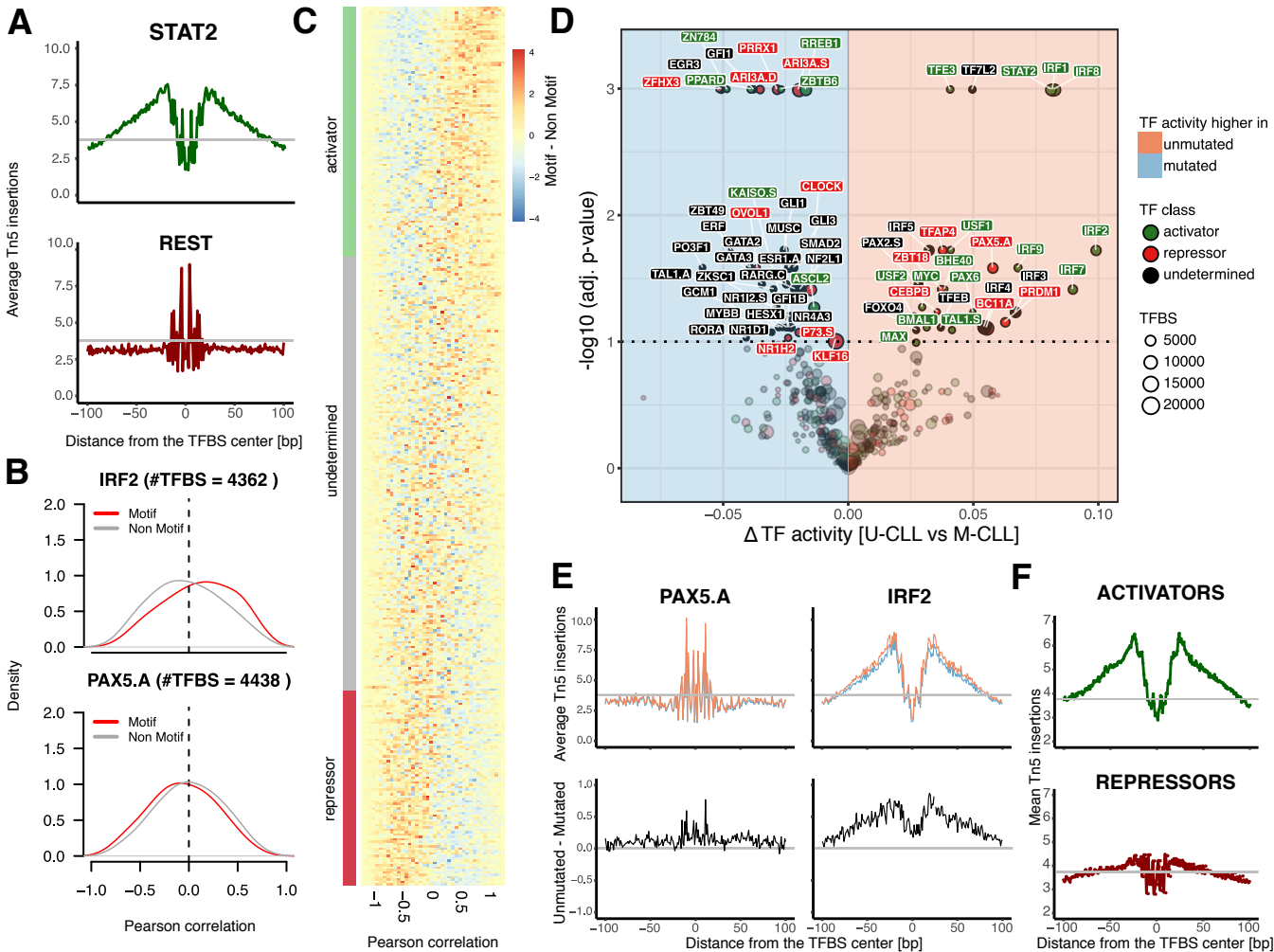
Significance 5% FDR: ○ 1 ● 2 ● 3 ● 4

G

540 **Figure 3.** *Clustering of TFs based on the similarity of their PWMs.* **(a)** Boxplot of PWM clusters
541 with at least 3 members as defined by RSAT for the differential TF activity between U-CLL and
542 M-CLL. In each cluster, the most negative and most positive TF is labeled. **(b)-(d):** RSAT
543 clustering output and tree for specific clusters. **(b)** Cluster 12 (POU family), the most distinct
544 cluster for M-CLL patients. **(c)** Cluster 17, which has two distinct subclusters, representing the
545 NFAT and NFkB family, respectively. **(d)** Cluster 40 (IRF family), the most distinct cluster for
546 U-CLL patients.

A**B****C****D**

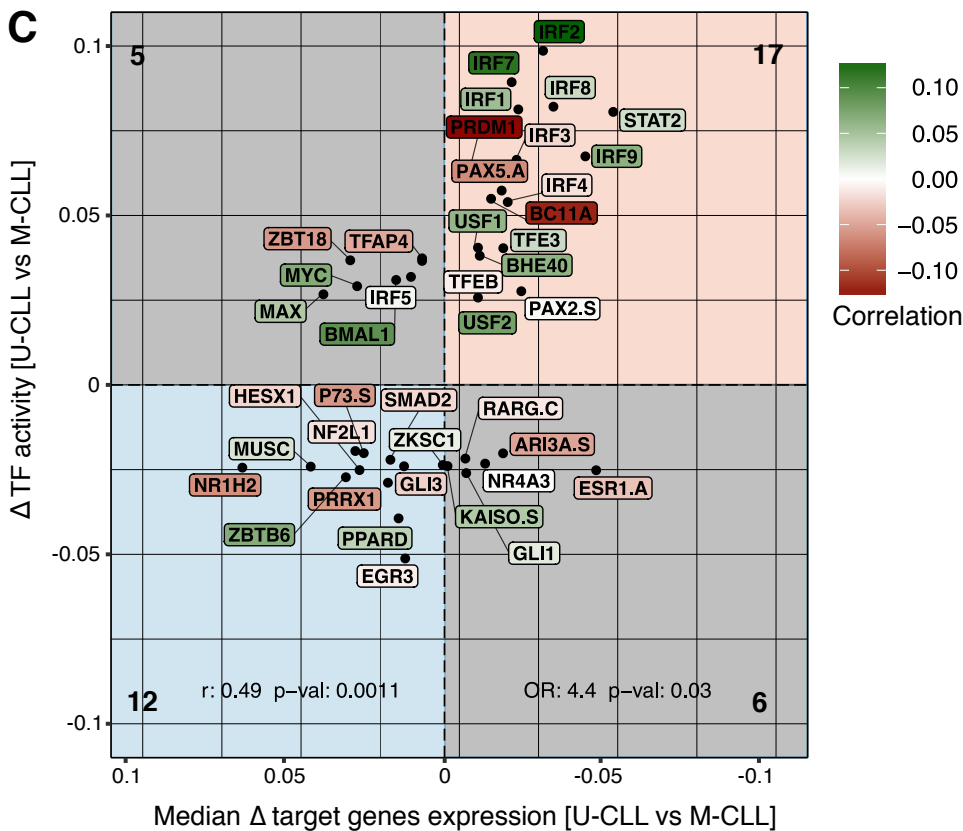
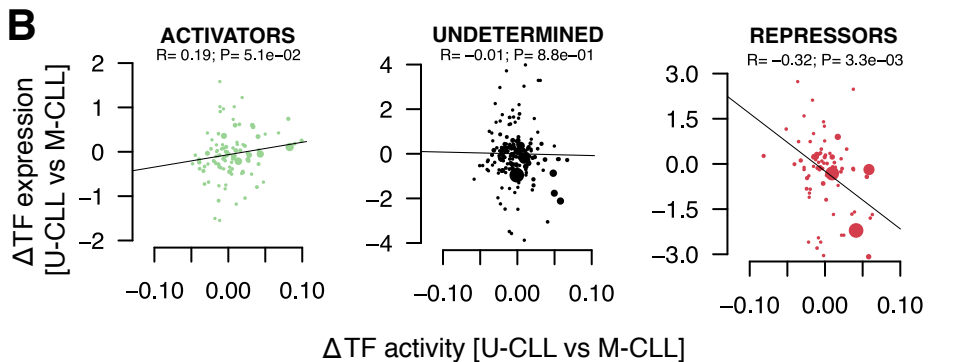
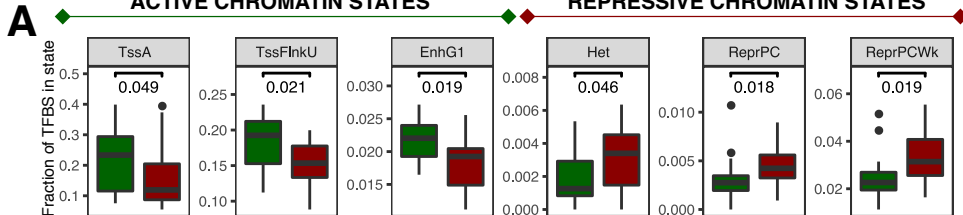
547 **Figure 4.** *Classification of TFs into activator or repressor based on RNA-Seq and ATAC-Seq*
548 *data.* **(a)** Exemplary footprints for a well-known activator (STAT2, top) and repressor (REST,
549 bottom). The x-axis depicts the distance in bp from the TFBS center, the y-axis denotes the
550 number of average Tn5 insertions, normalised to the library size and numbers of samples
551 across U-CLL and M-CLL. TFBS were predicted by *PWMscan* and only those overlapping with
552 open chromatin have been considered. The solid black line indicates the average insertion sites
553 within accessible chromatin. **(b)** Distributions of the Pearson correlations between TF
554 expression and ATAC-Seq signal at all putative TFBS (red line), and background distribution of
555 TFBS not containing the motif of interest (grey line) for two specific TFs, IRF2 (top) and PAX5
556 (bottom). **(c)** Summary of (b) across all TFs in the form of a heatmap showing the differences in
557 Pearson correlation of putative target peaks and background for each TF (i.e., subtracting the
558 black from the red line in (b)). Each TF is one row and is annotated as activator (green),
559 undetermined (grey) or repressor (red) as classified by *diffTF*. **(d)** Same as Fig. 2a, but with the
560 TFs labeled with their predicted role as activator (green), undetermined (black) or repressor
561 (red). See Fig. 2a for details. **(e)** Footprint analysis for an activator (IRF2, right) and repressor
562 (PAX5.A, left) as classified by *diffTF*. The top row shows the footprints separately for M-CLL and
563 U-CLL (blue and orange, respectively) based on the normalized number of Tn5 insertions, while
564 the bottom row highlights their differences (U-CLL - M-CLL) are shown. See (a) for axis
565 descriptions. **(f)** Summary footprint for all activators (top, green) and repressors (bottom, red).
566 See (a) for axis descriptions.



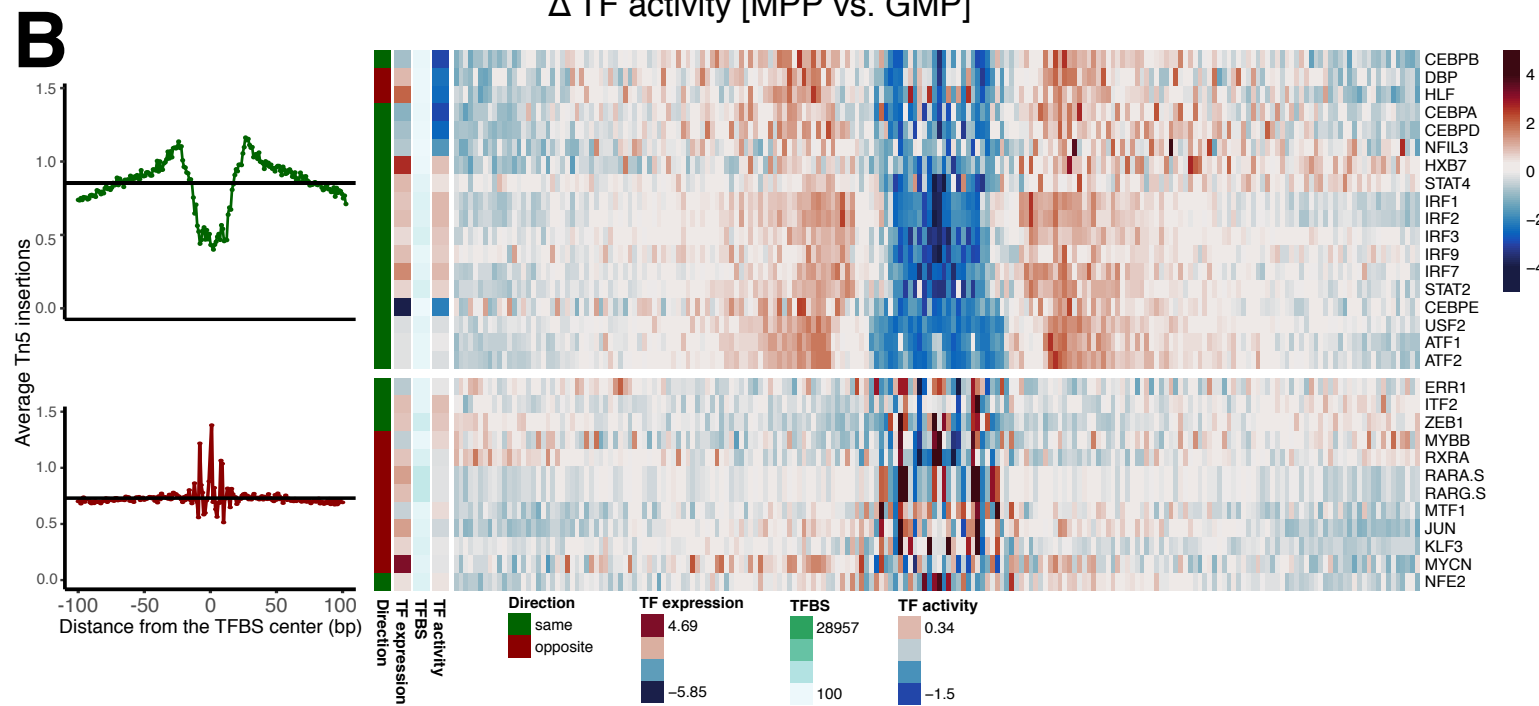
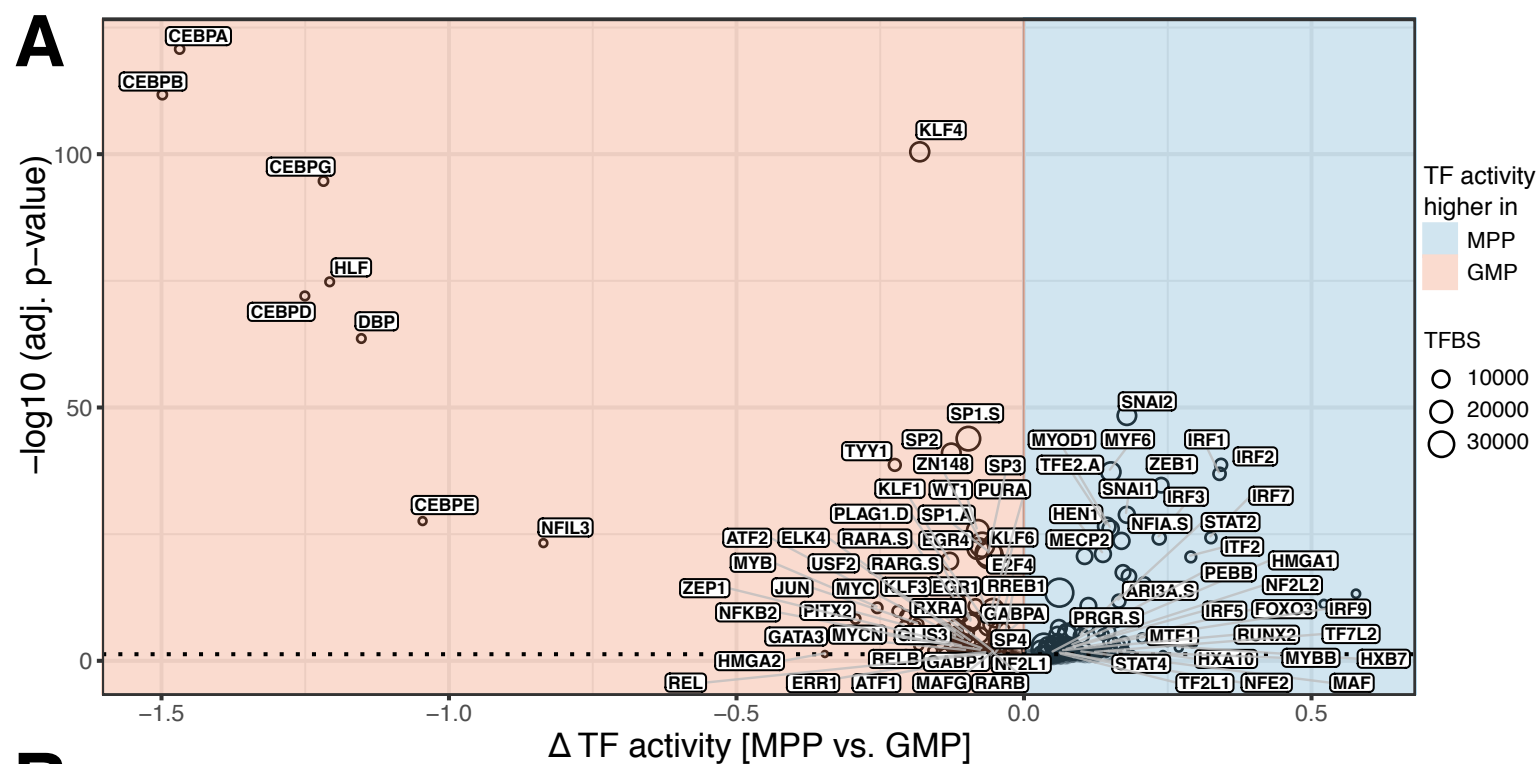
567 **Figure 5. Validations for the activator and repressor classification and downstream analyses. (a)**
568 Boxplots showing the fraction TFBS overlapping with specific chromatin states as defined by
569 *chromHMM* (Roadmap Epigenomics Consortium et al., 2015) for all activators (green) and
570 repressors (red) are shown. Only *chromHMM* states with significant differences (Wilcoxon test;
571 p-value < 0.05) between activators and repressors are displayed. **(b)** Pearson correlation of the
572 log₂ fold changes from RNA-Seq and differential TF activity for activators (green, left),
573 undetermined TFs (black, middle), and repressors (red, right). Only expressed TFs are shown.
574 **(c)** Correlation of differential TF activity and median of the differential target gene expression of
575 U-CLL against M-CLL. The x-axis shows the median target gene (TG) log₂ fold-change, the
576 y-axis denotes the differential TF activity. Each TF label is colored based on its
577 activator/repressor status (green/red) on a continuous scale (dark to light) based on the
578 correlation strength (see Fig. 4d). (OR=odds ratio, r=Pearson correlation coefficient).

ACTIVE CHROMATIN STATES

REPRESSIVE CHROMATIN STATES



579 **Figure 6.** *diffTF* recapitulates known TFs that drive the differentiation from MPP to GMP and
580 shows a similar activator/repressor cluster as in the CLL data. **(a)** Volcano plot of differential TF
581 activity between MPP (n=4) and GMP (n=4) cells. Due to the high number of significant TFs
582 only the most significant are labeled. The full list is available in Suppl. Table 4. **(b)** The footprints
583 are for TFs in two selected clusters that represent the activators and repressors, respectively
584 are shown as heatmap (right) and aggregate plots (left; see **Suppl. Fig. 17** for the full heatmap).
585 Only TFs that were significantly differentially active and all significantly differentially expressed
586 (adj. p-value < 0.05 for both) are displayed. Colors represent footprint strength, while white
587 denotes the value of the genomic background in the consensus peakset. Clusters were defined
588 using hierarchical clustering with the *ward.D2* method (clustering tree omitted for clarity)(see
589 **Suppl. Fig. 17**). For the cluster summary footprints at the left, we divided each footprint value by
590 the mean value of each cluster to highlight the differences in the surrounding chromatin
591 structure. The direction of TF expression and TF activity is in analogy to what is described in the
592 text. “Direction” denotes whether expression and TF activity have the same or opposite sign.



593 **REFERENCES**

- 594 Ambrosini, G., Groux, R., and Bucher, P. (2018). PWMScan: a fast tool for scanning entire
595 genomes with a position-specific weight matrix. *Bioinformatics* *34*, 2483–2484.
- 596 Arvaniti, E., Ntoufa, S., Papakonstantinou, N., Touloumenidou, T., Laoutaris, N.,
597 Anagnostopoulos, A., Lamnissou, K., Caligaris-Cappio, F., Stamatopoulos, K., Ghia, P., et al.
598 (2011). Toll-like receptor signaling pathway in chronic lymphocytic leukemia: distinct gene
599 expression profiles of potential pathogenic significance in specific subsets of patients.
600 *Haematologica* *96*, 1644–1652.
- 601 Baek, S., Goldstein, I., and Hager, G.L. (2017). Bivariate genomic footprinting detects changes
602 in transcription factor activity. *Cell Rep.* *19*, 1710–1722.
- 603 Baker, S.J., Ma'ayan, A., Lieu, Y.K., John, P., Reddy, M.V.R., Chen, E.Y., Duan, Q., Snoeck,
604 H.-W., and Reddy, E.P. (2014). B-myb is an essential regulator of hematopoietic stem cell and
605 myeloid progenitor cell development. *Proc Natl Acad Sci USA* *111*, 3122–3127.
- 606 Bonhoure, N., Bounova, G., Bernasconi, D., Praz, V., Lammers, F., Canella, D., Willis, I.M., Herr,
607 W., Hernandez, N., Delorenzi, M., et al. (2014). Quantifying ChIP-seq data: a spiking method
608 providing an internal reference for sample-to-sample normalization. *Genome Res.* *24*,
609 1157–1168.
- 610 Boorsma, A., Lu, X.-J., Zakrzewska, A., Klis, F.M., and Bussemaker, H.J. (2008). Inferring
611 condition-specific modulation of transcription factor activity in yeast through regulon-based
612 analysis of genomewide expression. *PLoS ONE* *3*, e3112.
- 613 Bussemaker, H.J., Li, H., and Siggia, E.D. (2001). Regulatory element detection using
614 correlation with expression. *Nat. Genet.* *27*, 167–171.
- 615 Chevrier, S., Emslie, D., Shi, W., Kratina, T., Wellard, C., Karnowski, A., Erikci, E., Smyth, G.K.,
616 Chowdhury, K., Tarlinton, D., et al. (2014). The BTB-ZF transcription factor Zbtb20 is driven by
617 Irf4 to promote plasma cell differentiation and longevity. *J. Exp. Med.* *211*, 827–840.
- 618 Chiorazzi, N., and Ferrarini, M. (2011). Cellular origin(s) of chronic lymphocytic leukemia:
619 cautionary notes and additional considerations and possibilities. *Blood* *117*, 1781–1791.
- 620 Coscia, M., Pantaleoni, F., Riganti, C., Vitale, C., Rigoni, M., Peola, S., Castella, B., Foglietta,
621 M., Griggio, V., Drandi, D., et al. (2011). IGHV unmutated CLL B cells are more prone to
622 spontaneous apoptosis and subject to environmental prosurvival signals than mutated CLL B
623 cells. *Leukemia* *25*, 828–837.
- 624 D'Annibale, S., Kim, J., Magliozzi, R., Low, T.Y., Mohammed, S., Heck, A.J.R., and
625 Guardavaccaro, D. (2014). Proteasome-dependent degradation of transcription factor activating
626 enhancer-binding protein 4 (TFAP4) controls mitotic division. *J. Biol. Chem.* *289*, 7730–7737.
- 627 El-Athman, R., and Relógio, A. (2018). Escaping circadian regulation: an emerging hallmark of
628 cancer? *Cell Syst.* *6*, 266–267.
- 629 Furman, R.R., Sharman, J.P., Coutre, S.E., Cheson, B.D., Pagel, J.M., Hillmen, P., Barrientos,
630 J.C., Zelenetz, A.D., Kipps, T.J., Flinn, I., et al. (2014). Idelalisib and rituximab in relapsed

- 631 chronic lymphocytic leukemia. *N. Engl. J. Med.* *370*, 997–1007.
- 632 Gascoyne, D.M., Long, E., Veiga-Fernandes, H., de Boer, J., Williams, O., Seddon, B., Coles,
633 M., Kioussis, D., and Brady, H.J.M. (2009). The basic leucine zipper transcription factor E4BP4
634 is essential for natural killer cell development. *Nat. Immunol.* *10*, 1118–1124.
- 635 Ghamlouch, H., Darwiche, W., Hodroge, A., Ouled-Haddou, H., Dupont, S., Singh, A.R.,
636 Guignant, C., Trudel, S., Royer, B., Gubler, B., et al. (2015). Factors involved in CLL
637 pathogenesis and cell survival are disrupted by differentiation of CLL B-cells into
638 antibody-secreting cells. *Oncotarget* *6*, 18484–18503.
- 639 Griffon, A., Barbier, Q., Dalino, J., van Helden, J., Spicuglia, S., and Ballester, B. (2015).
640 Integrative analysis of public ChIP-seq experiments reveals a complex multi-cell regulatory
641 landscape. *Nucleic Acids Res.* *43*, e27.
- 642 Grubert, F., Zaugg, J.B., Kasowski, M., Ursu, O., Spacek, D.V., Martin, A.R., Greenside, P.,
643 Srivas, R., Phanstiel, D.H., Pekowska, A., et al. (2015). Genetic control of chromatin states in
644 humans involves local and distal chromosomal interactions. *Cell* *162*, 1051–1065.
- 645 Han, H., Shim, H., Shin, D., Shim, J.E., Ko, Y., Shin, J., Kim, H., Cho, A., Kim, E., Lee, T., et al.
646 (2015). TRRUST: a reference database of human transcriptional regulatory interactions. *Sci.*
647 *Rep.* *5*, 11432.
- 648 Han, H., Cho, J.-W., Lee, S., Yun, A., Kim, H., Bae, D., Yang, S., Kim, C.Y., Lee, M., Kim, E., et
649 al. (2018). TRRUST v2: an expanded reference database of human and mouse transcriptional
650 regulatory interactions. *Nucleic Acids Res.* *46*, D380–D386.
- 651 Havelange, V., Pekarsky, Y., Nakamura, T., Palamarchuk, A., Alder, H., Rassenti, L., Kipps, T.,
652 and Croce, C.M. (2011). IRF4 mutations in chronic lymphocytic leukemia. *Blood* *118*,
653 2827–2829.
- 654 Heckman, C.A., Duan, H., Garcia, P.B., and Boxer, L.M. (2006). Oct transcription factors
655 mediate t(14;18) lymphoma cell survival by directly regulating bcl-2 expression. *Oncogene* *25*,
656 888–898.
- 657 Heinz, S., Benner, C., Spann, N., Bertolino, E., Lin, Y.C., Laslo, P., Cheng, J.X., Murre, C.,
658 Singh, H., and Glass, C.K. (2010). Simple combinations of lineage-determining transcription
659 factors prime cis-regulatory elements required for macrophage and B cell identities. *Mol. Cell*
660 *38*, 576–589.
- 661 Kasim, V., Xie, Y.-D., Wang, H.-M., Huang, C., Yan, X.-S., Nian, W.-Q., Zheng, X.-D., Miyagishi,
662 M., and Wu, S.-R. (2017). Transcription factor Yin Yang 2 is a novel regulator of the p53/p21
663 axis. *Oncotarget* *8*, 54694–54707.
- 664 Kauffman, E.C., Ricketts, C.J., Rais-Bahrami, S., Yang, Y., Merino, M.J., Bottaro, D.P.,
665 Srinivasan, R., and Linehan, W.M. (2014). Molecular genetics and cellular features of TFE3 and
666 TFEB fusion kidney cancers. *Nat. Rev. Urol.* *11*, 465–475.
- 667 Kern, D., Regl, G., Hofbauer, S.W., Altenhofer, P., Achatz, G., Dlugosz, A., Schnidar, H., Greil,
668 R., Hartmann, T.N., and Aberger, F. (2015). Hedgehog/GLI and PI3K signaling in the initiation

- 669 and maintenance of chronic lymphocytic leukemia. *Oncogene* *34*, 5341–5351.
- 670 Kikushige, Y., Ishikawa, F., Miyamoto, T., Shima, T., Urata, S., Yoshimoto, G., Mori, Y., Iino, T.,
671 Yamauchi, T., Eto, T., et al. (2011). Self-renewing hematopoietic stem cell is the primary target in
672 pathogenesis of human chronic lymphocytic leukemia. *Cancer Cell* *20*, 246–259.
- 673 Kim, R., Emi, M., and Tanabe, K. (2007). Cancer immunoediting from immune surveillance to
674 immune escape. *Immunology* *121*, 1–14.
- 675 Komatsu, M., Kurokawa, H., Waguri, S., Taguchi, K., Kobayashi, A., Ichimura, Y., Sou, Y.-S.,
676 Ueno, I., Sakamoto, A., Tong, K.I., et al. (2010). The selective autophagy substrate p62
677 activates the stress responsive transcription factor Nrf2 through inactivation of Keap1. *Nat. Cell*
678 *Biol.* *12*, 213–223.
- 679 Kreslavsky, T., Vilagos, B., Tagoh, H., Poliakova, D.K., Schwickert, T.A., Wöhner, M., Jaritz, M.,
680 Weiss, S., Taneja, R., Rossner, M.J., et al. (2017). Essential role for the transcription factor
681 Bhlhe41 in regulating the development, self-renewal and BCR repertoire of B-1a cells. *Nat.*
682 *Immunol.* *18*, 442–455.
- 683 Kulakovskiy, I.V., Medvedeva, Y.A., Schaefer, U., Kasianov, A.S., Vorontsov, I.E., Bajic, V.B.,
684 and Makeev, V.J. (2013). HOCOMOCO: a comprehensive collection of human transcription
685 factor binding sites models. *Nucleic Acids Res.* *41*, D195-202.
- 686 Landau, D.A., Tausch, E., Taylor-Weiner, A.N., Stewart, C., Reiter, J.G., Bahlo, J., Kluth, S.,
687 Bozic, I., Lawrence, M., Böttcher, S., et al. (2015). Mutations driving CLL and their evolution in
688 progression and relapse. *Nature* *526*, 525–530.
- 689 Landt, S.G., Marinov, G.K., Kundaje, A., Kheradpour, P., Pauli, F., Batzoglou, S., Bernstein,
690 B.E., Bickel, P., Brown, J.B., Cayting, P., et al. (2012). ChIP-seq guidelines and practices of the
691 ENCODE and modENCODE consortia. *Genome Res.* *22*, 1813–1831.
- 692 Li, Y.J., Sun, L., Shi, Y., Wang, G., Wang, X., Dunn, S.E., Iorio, C., Sreaton, R.A., and Spaner,
693 D.E. (2017). PPAR-delta promotes survival of chronic lymphocytic leukemia cells in energetically
694 unfavorable conditions. *Leukemia* *31*, 1905–1914.
- 695 Liu, L., Jin, G., and Zhou, X. (2015). Modeling the relationship of epigenetic modifications to
696 transcription factor binding. *Nucleic Acids Res.* *43*, 3873–3885.
- 697 Liu, Z., Lee, J., Krummey, S., Lu, W., Cai, H., and Lenardo, M.J. (2011). The kinase LRRK2 is a
698 regulator of the transcription factor NFAT that modulates the severity of inflammatory bowel
699 disease. *Nat. Immunol.* *12*, 1063–1070.
- 700 Medina-Rivera, A., Defrance, M., Sand, O., Herrmann, C., Castro-Mondragon, J.A., Delerce, J.,
701 Jaeger, S., Blanchet, C., Vincens, P., Caron, C., et al. (2015). RSAT 2015: regulatory sequence
702 analysis tools. *Nucleic Acids Res.* *43*, W50-6.
- 703 Mezger, A., Klemm, S., Mann, I., Brower, K., Mir, A., Bostick, M., Farmer, A., Fordyce, P.,
704 Linnarsson, S., and Greenleaf, W. (2018). High-throughput chromatin accessibility profiling at
705 single-cell resolution. *Nat. Commun.* *9*, 3647.
- 706 Minami, Y., Oishi, I., Endo, M., and Nishita, M. (2010). Ror-family receptor tyrosine kinases in

- 707 noncanonical Wnt signaling: their implications in developmental morphogenesis and human
708 diseases. *Dev. Dyn.* 239, 1–15.
- 709 Mittal, A.K., Chaturvedi, N.K., Rohlfesen, R.A., Gupta, P., Joshi, A.D., Hegde, G.V., Bociek, R.G.,
710 and Joshi, S.S. (2013). Role of CTLA4 in the proliferation and survival of chronic lymphocytic
711 leukemia. *PLoS ONE* 8, e70352.
- 712 Neu, K.E., and Wilson, P.C. (2016). Taking the broad view on B cell affinity maturation. *Immunity*
713 44, 518–520.
- 714 Oakes, C.C., Seifert, M., Assenov, Y., Gu, L., Przekopowicz, M., Ruppert, A.S., Wang, Q.,
715 Imbusch, C.D., Serva, A., Koser, S.D., et al. (2016). DNA methylation dynamics during B cell
716 maturation underlie a continuum of disease phenotypes in chronic lymphocytic leukemia. *Nat.*
717 *Genet.* 48, 253–264.
- 718 Rasmussen, K., Berest, I., Kessler, S., Nishimura, K., Simon-Carrasco, L., Vassilou, G.S.,
719 Pedersen, M.T., Christensen, J., Zaugg, J., and Helin, K. (2018). TET2 binding to enhancers
720 facilitates transcription factor recruitment in hematopoietic cells. *BioRxiv*.
- 721 Rendeiro, A.F., Schmidl, C., Strefford, J.C., Walewska, R., Davis, Z., Farlik, M., Oscier, D., and
722 Bock, C. (2016). Chromatin accessibility maps of chronic lymphocytic leukaemia identify
723 subtype-specific epigenome signatures and transcription regulatory networks. *Nat. Commun.* 7,
724 11938.
- 725 Roadmap Epigenomics Consortium, Kundaje, A., Meuleman, W., Ernst, J., Bilenky, M., Yen, A.,
726 Heravi-Moussavi, A., Kheradpour, P., Zhang, Z., Wang, J., et al. (2015). Integrative analysis of
727 111 reference human epigenomes. *Nature* 518, 317–330.
- 728 Sands, W.A., Copland, M., and Wheadon, H. (2013). Targeting self-renewal pathways in myeloid
729 malignancies. *Cell Commun. Signal.* 11, 33.
- 730 Schep, A.N., Wu, B., Buenrostro, J.D., and Greenleaf, W.J. (2017). chromVAR: inferring
731 transcription-factor-associated accessibility from single-cell epigenomic data. *Nat. Methods* 14,
732 975–978.
- 733 Slager, S.L., Caporaso, N.E., de Sanjose, S., and Goldin, L.R. (2013). Genetic susceptibility to
734 chronic lymphocytic leukemia. *Semin. Hematol.* 50, 296–302.
- 735 Sun, J., Ramos, A., Chapman, B., Johnnidis, J.B., Le, L., Ho, Y.-J., Klein, A., Hofmann, O., and
736 Camargo, F.D. (2014). Clonal dynamics of native haematopoiesis. *Nature* 514, 322–327.
- 737 Teng, M.W.L., Swann, J.B., Koebel, C.M., Schreiber, R.D., and Smyth, M.J. (2008).
738 Immune-mediated dormancy: an equilibrium with cancer. *J. Leukoc. Biol.* 84, 988–993.
- 739 Whyte, W.A., Orlando, D.A., Hnisz, D., Abraham, B.J., Lin, C.Y., Kagey, M.H., Rahl, P.B., Lee,
740 T.I., and Young, R.A. (2013). Master transcription factors and mediator establish
741 super-enhancers at key cell identity genes. *Cell* 153, 307–319.
- 742 Yanai, H., Negishi, H., and Taniguchi, T. (2012). The IRF family of transcription factors:
743 Inception, impact and implications in oncogenesis. *Oncoimmunology* 1, 1376–1386.

- 744 Yasuda, T., Hayakawa, F., Kurahashi, S., Sugimoto, K., Minami, Y., Tomita, A., and Naoe, T.
745 (2012). B cell receptor-ERK1/2 signal cancels PAX5-dependent repression of BLIMP1 through
746 PAX5 phosphorylation: a mechanism of antigen-triggering plasma cell differentiation. *J.*
747 *Immunol.* 188, 6127–6134.
- 748 Yeomans, A., Thirdborough, S.M., Valle-Argos, B., Linley, A., Krysov, S., Hidalgo, M.S.,
749 Leonard, E., Ishfaq, M., Wagner, S.D., Willis, A.E., et al. (2016). Engagement of the B-cell
750 receptor of chronic lymphocytic leukemia cells drives global and MYC-specific mRNA
751 translation. *Blood* 127, 449–457.
- 752 Zhou, Y., Li, Y.-S., Bandi, S.R., Tang, L., Shinton, S.A., Hayakawa, K., and Hardy, R.R. (2015).
753 Lin28b promotes fetal B lymphopoiesis through the transcription factor Arid3a. *J. Exp. Med.* 212,
754 569–580.