# Supplementary Notes

From subsection: **An atlas of immune cells in resting and stimulated states**

We observed high technical reproducibility (ATAC and RNA mean Pearson's R value = 0.89 and 0.80, respectively) and biological reproducibility (ATAC and RNA, mean Pearson's R value = 0.85 and 0.77, respectively) across replicates. We further confirmed the quality of our data by analyzing the enrichment of ATAC-seq reads mapping to transcription start sites (TSSs), and the expression of cell type-specific genes (Supplementary Table 1). As expected, we observed strong enrichment of reads at TSSs genome-wide and at promoters of cell type-specific genes such as *CD8A* in CD8+ T cells (Fig. 1d, Supplementary Fig. 2a, b).

From subsection: **Identifying immune memory-associated accessible regions**

Memory CD4+ T effectors had 2,275 peaks that exhibited increased accessibility compared to naive CD4+ T effector cells and only 130 peaks showed significantly decreased accessibility. In contrast to the increase in accessibility during maturation observed in lymphocytes, immature NK cells transitioning to memory NK cells gained accessibility at 1,211 sites while losing 2,526 accessible sites.

Characterizing peaks that correspond to lineage-specific and shared memory components in more detail, we observed signatures characteristic of enhancers in multiple blood-related cell types (Supplementary Fig. 3). Specifically, candidate binding sites of TFs with known roles in regulating memory T cell formation, such as RUNX3[1], BCL6[2,3] and NFKB1[4], showed strong enrichment in peaks that exhibit increased accessibility in T and B memory cells.

From subsection: **Stimulation leads to large-scale chromatin changes**

Connecting chromatin changes to gene expression, we found significant correlation between promoter accessibility (defined as a 5kb window around the start of a gene) and gene expression in resting state samples (R = 0.4) and stimulated samples (R = 0.37) (Supplementary Fig. 8b). Furthermore, we observed significant increases in accessibility in promoter regions of genes with the largest increases in expression upon stimulation (Supplementary Fig. 8c). Overall, these results illustrate the global chromatin and transcriptional changes of immune cells upon stimulation.

We next sought to identify transcription factors that may drive cell type and stimulation-responsive elements. We investigated variation in accessibility at position weight matrix (PWM)-predicted TF binding regions across cell types and conditions (Supplementary Fig. 5d). For example, the SPI1 motif is most enriched in B cells, DCs, and monocytes, consistent with gene expression data (Supplementary Fig. 5e). The corresponding TF for the SPI1 motif is integral to both myeloid and lymphoid B cell development[5].

In contrast, we found that the BATF motif (or perhaps another TF from the AP-1 family, which can have similar PWMs) was consistently enriched in accessibility regions across stimulated samples compared to their corresponding unstimulated state (Supplementary Fig. 5d). This suggests a shared effect of BATF and/or related transcription factors on chromatin regulation in stimulated samples across cell lineages, which was previously identified in stimulated CD4+ T cells[6]. Additionally, the putative activity of BATF correlates with upregulation of *BATF* expression (Supplementary Fig. 5f) and the expression of several classes of previously identified BATF-target genes (Supplementary Fig. 5g)[7]. Thus, our analysis identified large-scale genome-wide changes in chromatin accessibility and gene expression upon stimulation in B and T cells putatively attributable to specific sets of TFs.

From section: **Discussion**

Interestingly, rs6927172 was not detected as an eQTL in GTEx v7, likely because tissues represent cell mixtures with generally low proportions of immune cells and even lower proportions in activated states. However, Wu et al. demonstrated that the disruption of an 11bp region harboring rs6927172 significantly decreased gene expression of *TNFAIP3* in stimulated HEK293 T cells[8], suggesting this variant drives *TNFAIP3* expression.

57
58 Upon stimulation, Coornaert et al.[9] found that *TNFAIP3* expression first decreases and then
59 reappears, suggesting that the initial removal of A20 (encoded by *TNFAIP3)* is essential for optimal
60 NFKB1 activation[9]. While the mechanism that leads to the opening of the region containing rs6927172
61 and the suppression of A20 is unclear, we propose a model in which A20 is first down-regulated by
62 other factors, allowing activation of NFKB1 (Supplementary Fig. 14). Next, NFKB1 binds to the region
63 containing rs6927162, resulting in the reappearance of A20 expression. Such a hypothesis is
64 supported by studies focused on the regulation of NFKB1[10]. Thus, rs6927172 likely prevents the return
65 of A20 by disrupting the binding of NFKB1, which subsequently results in inappropriate NFKB1
66 signaling.
67

# Additional Methods

69
70 **Data collection**
71
72 *Fetal sample processing*
73 Human thymus was obtained from 18- to 22-gestational-week specimens under the guidelines of the
74 Committee on Human Research (UCSF IRB)–approved protocols from the Department of Obstetrics,
75 Gynecology and Reproductive Science, San Francisco General Hospital. Fetal samples were obtained
76 after legal, elective termination of pregnancy with written informed consent for fetal tissue donation to
77 biomedical research. Consent for tissue donation was obtained by clinical staff after the decision to
78 pursue termination was reached by patients. Personal Health Information and Medical Record
79 Identifiers/access is at no point available to researchers, and no such information is associated with
80 tissue samples at any point. Tissue was washed and cut into small pieces using scissors. Thymocytes
81 were extracted by mashing tissue pieces gently using the back of a sterile syringe. To extract TECs,
82 remaining tissue pieces were digested for 30 min at 37°C using medium containing 100 µg ml$^{-1}$
83 DNase I (Roche, Switzerland) and 100 µg ml$^{-1}$ Liberase TM (Sigma-Aldrich, MO, USA) in RPMI.
84 Fragments were triturated through a 5-ml pipette every 6 min to mechanically aid digestion. At 30 min,
85 tubes were spun briefly to pellet undigested fragments and the supernatant was discarded. Fresh
86 digestion medium was added to remaining fragments and the digestion was repeated using a glass
87 Pasteur pipette for trituration. Supernatant from this second round of digestion was also discarded. A
88 third round of enzymatic digestion was performed using digestion medium supplemented with trypsin-
89 EDTA for a final concentration of 0.05%. Remaining thymic fragments were digested for another 30
90 min or until a single cell suspension was obtained. The cells were moved to cold MACS buffer (0.5%
91 BSA, 2 mM EDTA in PBS) to stop the enzymatic digestion. Following digestion, TECs were enriched
92 by density centrifugation over a three-layer Percoll gradient with specific gravities of 1.115, 1.065 and
93 1.0. Stromal cells isolated from the Percoll-light fraction (between the 1.065 and 1.0 layers) were
94 washed in MACS buffer. Samples were sorted on FACS Aria flow cytometer (BD Biosciences, CA,
95 USA) up to >95% purity. Sorted cells were washed once in PBS, cryopreserved in Bambanker freezing
96 media (LYMPHOTEC Inc, Japan) for ATAC experiments and in TriReagent (Sigma-Aldrich, MO, USA)
97 for RNA experiments. Cells frozen in Bambanker freezing media were stored in liquid nitrogen until
98 ready to use. Cells frozen in TriReagent were stored at -80°C until further use.
99
100 *NFKB1 ChIP-seq from heterozygous donors*
101 Isolated CD4$^+$ T cells were stimulated for 24 hours with anti-human CD3/CD28 dynabeads (Thermo
102 Fisher Scientific, MA, USA) at a 1:1 cell to bead ratio and 50 unit/ml of human IL-2 (UCSF Pharmacy).
103 The cells were harvested, fixed with 1% formaldehyde (Thermo Fisher Scientific, MA, USA) for 10 min,
104 washed twice with cold PBS and frozen at -80C. The chromatin was sonicated to generate fragments
105 of 200-500bp in length, followed by incubation with rabbit polyclonal p50 (#3035, Cell Signaling,
106 Danvers, MA, USA) and p65 (ab16502, Abcam, Cambridge, UK) antibodies overnight. Precipitated
107 chromatin was washed, de-crosslinked and DNA extraction was carried out using phenol-chloroform
108 (Sigma). ChIP DNA was prepared for high throughput sequencing using Accel-NGS 2S Plus DNA
109 library kit (Swift Biosciences) as per manufacturer's protocol. DNA libraries were sequenced on an
110 Illumina Hiseq4000 with a paired-end 75bp run (CAT, UCSF).
111
112 **Collection of publicly available data**
113

*Progenitor RNA-seq and ATAC-seq data from GEO*
115 ATAC-seq and RNA-seq of hematopoiesis progenitors and several differentiated cell types were
116 downloaded[11], and processed through the respective ATAC-seq and RNA-seq pipeline described
117 below. Only data from healthy controls was included throughout this study.
118
119 *ATAC-seq data from ENCODE*
120 To serve as a negative control for GWAS enrichment analyses, we collected data from ATAC-seq
121 samples from tissues with low proportions of immune cells. We used the ENCODE data portal to
122 download all available raw fastq ATAC-seq files from the calf muscle and breast epithelium human
123 tissues. All data were processed with the same ATAC-seq data processing pipeline described below.
124
125 *Obtaining GWAS summary statistics*
126 We downloaded the full set of GWAS summary statistics of 13 complex traits from 9 autoimmune traits
127 and 4 primarily non-autoimmune and thus negative control traits (Sunburn, Alzheimer's disease, Type
128 2 diabetes, and Schizophrenia) that have been previously aggregated[12]. For analyses that relied on
129 fine-mapped disease-associated variants, we downloaded the list of PIC variants. These represent
130 individual GWAS regions that have been fine-mapped with a previously described statistical method
131 [13].
132
133 *Obtaining Blood eQTL summary statistics*
134 From the GTEx data portal we downloaded v7 eQTL estimates for all SNP-gene pairs tested with
135 whole blood gene expression.
136
137 **Data analysis**
138
139 *ChIP-seq*
140 We aligned ChIP-seq reads using bowtie2 version 2.2.9 with default parameters and a maximum
141 paired-end insert distance of 2kbp. The bowtie2 index was constructed with the default parameters for
142 the hg19 reference genome. We filtered out reads that mapped to chrM and used samtools version 1.4
143 to filter out reads with MAPQ < 30 and with the flags '-F 1804' and '-f 2'. Additionally, duplicate reads
144 were discarded using picard version 1.134 (http://broadinstitute.github.io/picard/). TF-bound peaks
145 were identified with MACS2 version 2.1.1 under default parameters and '--nomodel --nolambda --
146 keep-dup all --call-summits'. Additionally, the appropriate input-DNA background was set with the '--
147 control' parameter. A consensus set of peaks was defined by merging overlapping (1bp or more)
148 peaks identified in at least two samples across all samples. We then used the 'get_count' function
149 from the nucleoATAC python package to count the number of fragments within the consensus peak
150 set across all samples[14]. We used samtools version 1.4 'mpileup' to count reads aligning to
151 rs6927172.
152
153 *Exploratory analysis*
154 We used tSNE, PCA, and k-means clustering to explore trends in gene expression and chromatin
155 accessibility variation. As input to these methods, we used the read count matrices corrected for
156 sample quality (with TSS enrichment as a proxy) and batch effects (with donor as a proxy). After
157 using trimmed mean of M-values (TMM) to estimate scaling factors[15] and applying the voom
158 transformation[16], we used the 'removeBatchEffect' function from limma to regress out batch and
159 sample quality effects [17]. Aside from the removal of these effects, normalized counts are equivalent to
160 the addition of a 0.5 pseudocount to the fragment count and a log2 transformation of the fragment
161 counts per million (CPM). For analyses that included previously published samples, we did not remove
162 batch effects, because these batches do not have overlapping cell types. However, batch appeared to
163 have a minimal effect on sample clusters and these analyses were exploratory in nature. The package
164 Rtsne version 0.13 (https://github.com/jkrijthe/Rtsne) was used for tSNE analysis with default
165 parameters, unless there were too few samples in which case the perplexity was set to 10.
166 Additionally, we performed k-means clustering (with k set to the total number of cell-type by condition
167 pairs) and then estimated the accuracy of ATAC-seq and RNA-seq unsupervised clustering by
168 computing the HA-adjusted RAND index[18] considering the known cell-type/condition pairs as the
169 ground truth clustering. We repeated the clustering 100 times to estimate the average HA-adjusted
170 RAND index. When comparing HA-adjusted RAND index values between the RNA-seq and ATAC-seq
171 samples, we used the intersection of samples.

172
173 *Correlation between promoter accessibility and gene expression*
174 For each sample, we counted the number of filtered ATAC-seq reads that aligned to 5kb promoter
175 regions based on annotations of unique protein coding genes from gencode v25. Once again, we used
176 trimmed mean of M-values (TMM) to estimate scaling factors and applied the voom transformation to
177 compute log2(CPM) counts following the addition of a 0.5 pseudocount. We merged samples from the
178 same cell type and condition across donors by averaging the log2(CPM) accessibility values for each
179 promoter region. Finally, we quantile normalized accessibility values to a standard normal distribution
180 along with the processed gene expression values. We reported Pearson's R correlation values to
181 assess the relationship between promoter accessibility and gene expression. Values presented in
182 Supplementary Fig. 8 are from all samples, however we report condition and cell type-specific
183 correlation values in Supplementary Table 1.
184
185 *Enrichment analysis of differentially accessible regions*
186 We used Fisher's exact tests implemented in the LOLA tool[19] to quantify enrichment of sets of
187 differentially accessible peaks in different lineages in comparison to a universe set of all peak regions
188 that were shared between comparisons. As a catalogue for potential enrichment we considered a
189 collection comprising peaks from the CODEX database[20], ENCODE TFBS, ENCODE chromatin state
190 segmentations and candidate binding sites for motifs in the JASPAR database[21] determined by the
191 motifmatchr R package (https://github.com/GreenleafLab/motifmatchr). The rank represents the rank
192 of each dataset for a given peak set. Max rank means, the largest (i.e. worst) rank among the following
193 scores from a Fisher's exact test: odd-ratio, p-value, support. The white squares represent non-
194 significant enrichments (q-value >= 0.01) or enrichments that could not be computed, because there
195 was no overlap.
196
197 *Enrichment analysis of differentially expressed genes*
198 For each subset we identified significantly differentially expressed genes with a *q*-value less than 0.01
199 and absolute log2FC greater than 1. We used g:Profiler to identify pathways that were significantly
200 enriched for stimulation-associated genes[22] with an ordered query based on a ranking of differential
201 expression q-values and Bonferroni p-value correction. Top enriched pathways per cell subset are
202 listed in Supplementary Table 1 and Supplementary Fig. 8 displays a heatmap of these results.
203
204 *Variance decomposition*
205 We were interested in decomposing the total variance of chromatin accessibility across samples into
206 variance components attributable to specific factors. Therefore, we fitted a random effects model:
207
208 $$a_{ij} = \kappa_{i,s(j)} + \beta_{i,l(j)} + \gamma_{i,c(j)} + \zeta_{i,s(j),l(j)} + \eta_{i,s(j),c(j)} + \delta_{i,d(j)} + \lambda_i t_j + \epsilon_{ij},$$
209
210 where chromatin accessibility ($a$) at peak $i$ for a sample $j$ is a function of the effects of the stimulation
211 condition ($\kappa$), lineage ($\beta$), cell type ($\gamma$), lineage/stimulation interaction ($\zeta$), cell/stimulation interaction
212 ($\eta$), donor ($\delta$), TSS enrichment ($t$) of ATAC-seq reads ($\lambda$), and the residual error ($\epsilon$). For notational
213 convenience, we define a function for each feature in the model that looks up sample-specific
214 information, i.e., $s(j)$ represents the stimulation condition associated with sample $j$. We represented
215 accessibility with the log2(CPM) ATAC-seq read counts (with the addition of a pseudocount of 0.5) at
216 consensus peaks across samples, which were normalized for read depth with TMM normalization.
217 Additionally, we scaled accessibility at each peak across samples to have mean=0 and variance=1.
218 We included the effects of $d$ and $t$ (as a proxy for sample quality), to control for their effects, since the
219 other parameters are our primary interest. Across all peaks, we modelled the distribution of effects:
220
221 $$(\kappa, \beta, \gamma, \zeta, \eta, \delta, \lambda, \epsilon) \sim MVN\big(0, diag(\sigma_s^2, \sigma_l^2, \sigma_c^2, \sigma_{sl}^2, \sigma_{sc}^2, \sigma_d^2, \sigma_t^2, \sigma_\epsilon^2)\big).$$
222
223 We used a maximum likelihood approach to jointly estimate the $\sigma^2$ parameters for each factor. To
224 obtain robust estimates we found it beneficial to pool peaks. We found pooling 100 peaks represented
225 a good compromise between computational cost and statistical robustness. To assess uncertainty of
226 variance estimates, we repeated the analysis on 100 sets of 100 randomly selected peaks with
227 replacement. The total biological variance explained (TBVE) by the factors of interest is,
228
229 $$TBVE = \sigma_s^2 + \sigma_l^2 + \sigma_c^2 + \sigma_{sl}^2 + \sigma_{sc}^2.$$

4

230
231 Therefore, the proportion of biological variance explained (PBVE) contributed by a factor is the
232 variance estimate $\sigma^2$ for that factor divided by $TBVE$. We listed the median value across all bootstrap
233 replicates. For results reported we limited our analysis to cell types from the four donors with the most
234 cell samples collected and excluded cell types with fewer than three biological replicates.
235
236 *Visualizing TF ATAC-seq footprints*
237 We aggregated ATAC-seq insertion counts around candidate binding sites for motifs in the JASPAR
238 database[21] determined by the motifmatchr R package (https://github.com/GreenleafLab/motifmatchr)
239 using transcription factor footprinting methods previously described[23].
240
241 *Boxplot visualizations*
242 Unless otherwise mentioned, all boxplot visualizations represent the median, two hinges (25[th] and 75[th]
243 percentile) and whiskers. Whiskers show a line from the hinge to 1.5 * the difference between the first
244 and third quartile. Points that extend beyond the whiskers are displayed individually.
245
246 *TF position weight matrix (PWM) motif analyses*
247 For determining PWMs enriched in open chromatin regions we used chromVAR version 1.0.1 with
248 default parameters on read counts within consensus peaks of samples merged by donor[24]. Following
249 identification of condition-associated TFs with chromVAR we wanted to examine the effect of a few of
250 these TFs on allele-specific chromatin accessibility. We used the PWM of a TF of interest to predict
251 the binding affinity of a 41 bp genomic region centered on the heterozygous site. The binding affinity or
252 match score was computed using the 'motifmatchr' R package
253 (https://github.com/GreenleafLab/motifmatchr), which is a wrapper for the MOODS motif matching
254 suite[25]. The relative binding score was determined by subtracting the binding affinity match score of
255 the alternative allele from that of the reference allele. As a threshold for presence or absence of motif
256 matching we used a p value cutoff of $3 \times 10^{-3}$. In this way we grouped heterozygous sites into three
257 groups: predicted TF affinity for the reference (relative match > 1), alternative (relative match < -1) or
258 no preference (absolute value of the relative match < 0.01).
259
260 *Peak clustering*
261 To test whether disease heritability was enriched within stimulation-specific chromatin accessible from
262 B and T cell lineages we used a supervised peak clustering approach. First, we scaled the matrix of
263 ATAC-seq read counts per sample (indicated by $j$) across all consensus peaks (indicated by $i$) to
264 values between 0 and 1 with

265
$$x'_{i,j} = \frac{x_{i,j} - \min(x_{i,\cdot})}{\max(x_{i,\cdot}) - \min(x_{i,\cdot})},$$

266
267 where $x'$ represents the scaled matrix. Peaks in a sample that were in the top decile were
268 automatically set to 1 to represent the fully accessible state. Per peak we computed the median scaled
269 accessibility across samples from the same broad cell type and condition.
270
271 Our goal was to identify peaks that express a specific accessibility profile. We defined a profile of
272 interest with a vector of length equal to the number of merged lineage and condition samples with
273 values of either 0 or 1 corresponding to closed or open chromatin accessibility. We consider 11
274 profiles of interest (Supplementary Fig. 10d). To identify peaks with a similar profile we computed the
275 average Euclidean distance between each of the ideal accessibility profiles and each peak. Peaks
276 more similar to an ideal peak profile should have smaller distances to the peak profile. Additionally,
277 when computing the distance, we incorporated a weight per sample to influence the importance of
278 matching accessibility in different merged samples. This was important to find peaks that were
279 accessible in resting samples (weight of 1), while allowing for the possibility that the peak was
280 accessible in the same lineage but stimulated samples (weight of 0).
281
282 To determine a distance cutoff of a peak cluster, we permuted the peak accessibility values for each
283 sample and computed a null distribution of distances for each lineage and peak cluster type. We used
284 a peak distance threshold resulting in fewer than 5% false positives. Finally, peaks passing this peak
285 distance threshold were assigned to a single profile of interest based on the minimum distance, thus
286 forming disjoint sets of accessible regions.

5

# References

1. Wang, D. *et al.* The Transcription Factor Runx3 Establishes Chromatin Accessibility of cis-Regulatory Landscapes that Drive Memory Cytotoxic T Lymphocyte Formation. *Immunity* **48**, 659-674 e6 (2018).
2. Ichii, H. *et al.* Role for Bcl-6 in the generation and maintenance of memory CD8+ T cells. *Nat Immunol* **3**, 558-63 (2002).
3. Fukuda, T. *et al.* Disruption of the Bcl6 gene results in an impaired germinal center formation. *J Exp Med* **186**, 439-48 (1997).
4. Lai, W. *et al.* Transcriptional control of rapid recall by memory CD4 T cells. *J Immunol* **187**, 133-40 (2011).
5. Lloberas, J., Soler, C. & Celada, A. The key role of PU.1/SPI-1 in B cells, myeloid cells and macrophages. *Immunol Today* **20**, 184-9 (1999).
6. Gate, R.E. *et al.* Genetic determinants of co-accessible chromatin regions in activated T cells across humans. *Nat Genet* (2018).
7. Kurachi, M. *et al.* The transcription factor BATF operates as an essential differentiation checkpoint in early effector CD8+ T cells. *Nat Immunol* **15**, 373-83 (2014).
8. Wu, J. *et al.* CRISPR/cas9 mediated knockout of an intergenic variant rs6927172 identified IL-20RA as a new risk gene for multiple autoimmune diseases. *Genes Immun* (2018).
9. Coornaert, B. *et al.* T cell antigen receptor stimulation induces MALT1 paracaspase-mediated cleavage of the NF-kappaB inhibitor A20. *Nat Immunol* **9**, 263-71 (2008).
10. Housley, W.J. *et al.* Genetic variants associated with autoimmunity drive NFkappaB signaling and responses to inflammatory stimuli. *Sci Transl Med* **7**, 291ra93 (2015).
11. Corces, M.R. *et al.* Lineage-specific and single-cell chromatin accessibility charts human hematopoiesis and leukemia evolution. *Nat Genet* **48**, 1193-203 (2016).
12. Reshef, Y.A. *et al.* Detecting genome-wide directional effects of transcription factor binding on polygenic disease risk. *Nat Genet* **50**, 1483-1493 (2018).
13. Farh, K.K. *et al.* Genetic and epigenetic fine mapping of causal autoimmune disease variants. *Nature* **518**, 337-43 (2015).
14. Schep, A.N. *et al.* Structured nucleosome fingerprints enable high-resolution mapping of chromatin architecture within regulatory regions. *Genome Res* **25**, 1757-70 (2015).
15. Robinson, M.D. & Oshlack, A. A scaling normalization method for differential expression analysis of RNA-seq data. *Genome Biol* **11**, R25 (2010).
16. Law, C.W., Chen, Y., Shi, W. & Smyth, G.K. voom: Precision weights unlock linear model analysis tools for RNA-seq read counts. *Genome Biol* **15**, R29 (2014).
17. Ritchie, M.E. *et al.* limma powers differential expression analyses for RNA-sequencing and microarray studies. *Nucleic Acids Res* **43**, e47 (2015).
18. Arabie, L.H.P. Comparing partitions. *Journal of Classification* **2**, 193-218 (1985).
19. Sheffield, N.C. & Bock, C. LOLA: enrichment analysis for genomic region sets and regulatory elements in R and Bioconductor. *Bioinformatics* **32**, 587-9 (2016).
20. Sanchez-Castillo, M. *et al.* CODEX: a next-generation sequencing experiment database for the haematopoietic and embryonic stem cell communities. *Nucleic Acids Res* **43**, D1117-23 (2015).
21. Khan, A. *et al.* JASPAR 2018: update of the open-access database of transcription factor binding profiles and its web framework. *Nucleic Acids Res* **46**, D1284 (2018).
22. Reimand, J. *et al.* g:Profiler-a web server for functional interpretation of gene lists (2016 update). *Nucleic Acids Res* **44**, W83-9 (2016).
23. Corces, M.R. *et al.* The chromatin accessibility landscape of primary human cancers. *Science* **362**(2018).
24. Schep, A.N., Wu, B., Buenrostro, J.D. & Greenleaf, W.J. chromVAR: inferring transcription-factor-associated accessibility from single-cell epigenomic data. *Nat Methods* **14**, 975-978 (2017).
25. Korhonen, J.H., Palin, K., Taipale, J. & Ukkonen, E. Fast motif matching revisited: high-order PWMs, SNPs and indels. *Bioinformatics* **33**, 514-521 (2017).