

1 **An experimentally validated network of nine haematopoietic transcription**
2 **factors reveals mechanisms of cell state stability**

3 Judith Schütte^{1,2}, Huange Wang¹, Stella Antoniou^{3,*}, Andrew Jarratt^{3,4,*}, Nicola K. Wilson¹,
4 Joey Riepsaame³, Fernando J. Calero-Nieto¹, Victoria Moignard¹, Silvia Basilico¹, Sarah J.
5 Kinston¹, Rebecca L. Hannah¹, Mun Chiang Chan³, Sylvia T. Nürnberg^{5,6}, Willem H.
6 Ouwehand⁵, Nicola Bonzanni^{7,\$}, Marella F.T.R. de Bruijn^{3,\$}, Berthold Göttgens^{1,\$}

7
8 ¹ Department of Haematology, Cambridge Institute for Medical Research, University of
9 Cambridge, Cambridge, United Kingdom and Wellcome Trust - Medical Research Council
10 Cambridge Stem Cell Institute, University of Cambridge, Cambridge, United Kingdom

11 ² Present address: Department of Haematology, University Hospital Essen, Essen, Germany

12 ³ MRC Molecular Haematology Unit, Weatherall Institute of Molecular Medicine, Radcliffe
13 Department of Medicine, University of Oxford, Oxford OX3 9DS, United Kingdom

14 ⁴ Present address: Division of Molecular Medicine, Walter and Eliza Hall Institute of Medical
15 Research, 1G Royal Parade, Parkville, VIC 3052, Australia

16 ⁵ Department of Haematology, University of Cambridge & NHS Blood and Transplant,
17 Cambridge, UK

18 ⁶ Present address: University of Pennsylvania, Perelman School of Medicine, Philadelphia PA
19 19104, USA

20 ⁷ IBIVU Centre for Integrative Bioinformatics, VU University Amsterdam, De Boelelaan 1081,
21 Amsterdam, The Netherlands and Netherlands Cancer Institute, Plesmanlaan 121, 1066 CX,
22 Amsterdam

23 * Contributed equally

24 \$ Authors for Correspondence

25 **Abstract**

26 Transcription factor (TF) networks determine cell type identity by establishing and maintaining
27 lineage-specific expression profiles, yet reconstruction of mammalian regulatory network
28 models has been hampered by a lack of comprehensive functional validation of regulatory
29 interactions. Here, we report comprehensive ChIP-Seq, transgenic and reporter gene
30 experimental data that have allowed us to construct an experimentally validated regulatory
31 network model for haematopoietic stem/progenitor cells (HSPCs). Model simulation coupled
32 with subsequent experimental validation using single cell expression profiling revealed
33 potential mechanisms for cell state stabilisation, and also how a leukemogenic TF fusion
34 protein perturbs key HSPC regulators. The approach presented here should help to improve
35 our understanding of both normal physiological and disease processes.

36 **Introduction**

37 Tight regulation of gene expression is essential for both the establishment and maintenance of
38 cellular phenotypes within metazoan organisms. The binding of transcription factor proteins
39 (TFs) to specific DNA sequence motifs represents the primary step of decoding genetic
40 information into specific gene expression patterns. TF binding sites (TFBSs) or motifs are
41 usually short (6-10 bp), and therefore found just by chance throughout the genome. Functional
42 TFBSs often occur as evolutionarily conserved clusters, which in the case of enhancers can act
43 over long distances, thus necessitating comprehensive analysis of entire gene loci to understand
44 the transcriptional control mechanisms acting at mammalian gene loci.

45 Given the complex regulatory circuitries that arise when control of multiple genes is
46 considered, transcriptional control is often represented in the form of gene regulatory networks
47 (GRNs), which carry most mechanistic information when constructed from detailed knowledge
48 on the TFs and the *cis*-regulatory elements with which they interact (1-6). Importantly,
49 regulatory network models can provide much more than a representation of existing
50 knowledge, because network simulations can reveal possible molecular mechanisms that
51 underlie highly complex biological processes. Boolean modelling approaches have been used
52 to reconstruct core regulatory networks in blood stem cells (7) and myeloid progenitors (8), but
53 neither of these studies took into account the underlying regulatory structure of the relevant
54 gene regulatory elements. Full gene-regulatory information has been used for an ordinary
55 differential equation-based model (9, 10), but was restricted to a small three-gene core circuit.
56 Large consortia efforts such as ImmGen and FANTOM5 have created comprehensive networks
57 of either regulatory elements or gene signatures important for multipotency and differentiation
58 (11, 12). Furthermore, studies looking at gene regulation circuitry in embryonic stem (ES) cells
59 have proposed regulatory networks important for ES cell identity (13, 14). While the
60 complexities of transcriptional control demand approaches such as network modelling, no

61 single experimental method can provide the complex biological data required for the
62 construction of accurate models. The previously mentioned studies focus their attention on one
63 specific aspect of network modelling and importantly did not combine network analysis with
64 comprehensive functional validation. Given that the key building blocks are gene regulatory
65 sequences and the TFs bound to them, essential information for network reconstruction
66 includes (i) comprehensive TF binding data, (ii) *in vivo* validation of the functionality of TF-
67 bound regions as bona fide regulatory elements, and (iii) molecular data on the functional
68 consequences of specific TF binding events (e.g. activation vs. repression). The regulatory
69 network model that we present in this study comprises all of the aforementioned components
70 and is accompanied by functional validation of model predictions.

71

72 **Results**

73 ***In vivo* validation of *cis*-elements as regulatory network nodes connecting 9 HSPC TFs**

74 For the reconstruction of a core GRN model for HSPCs, we focussed on nine major HSPC
75 regulators (ERG, FLI1, GATA2, GFI1B, LYL1, MEIS1, PU.1, RUNX1, TAL1), for which
76 genome-wide binding patterns in the murine multipotent progenitor cell line HPC7 have
77 previously been published (15). First, we searched the literature to summarise known *cis*-
78 regulatory regions for the nine TFs that possess haematopoietic activity in transgenic mouse
79 embryos, which recovered a total of 14 regions: *Erg*+85 (16), *Fli1*-15 (17), *Fli1*+12, *Gata2*-3
80 (18), *Gata2*+3 (= *Gata2*+9.5) (19), *Gfi1b*+13, *Gfi1b*+16, *Gfi1b*+17 (16, 20), *Lyl1* promoter
81 (21), *Spi1*-14 (22), *Runx1*+23 (23), *Tal1*-4 (24), *Tal1*+19 (25) and *Tal1*+40 (26).

82 To extend this partial knowledge of relevant gene regulatory sequences to a comprehensive
83 definition of how these nine TFs might cross-regulate each other, we made use of the genome-
84 wide binding data for the nine TFs (15) as well as information on acetylation of histone H3 at

85 lysine 27 (H3K27ac) (27) in the HPC7 blood progenitor cell line. Additional candidate gene
86 regulatory regions for all nine TFs were selected based on binding of at least three TFs and
87 H3K27ac, since it has been shown previously that transcriptionally active regions are
88 commonly bound by multiple TFs and display H3K27 acetylation (28). To assign putative
89 candidate regions to a given TF, they had to be located between its respective upstream and
90 downstream flanking genes, i.e. within the gene body itself or its 5' and 3' intergenic flanking
91 regions. The *Erg* gene locus for example contains five candidate *cis*-regulatory regions based
92 on these criteria, namely *Erg*+65, *Erg*+75, *Erg*+85, *Erg*+90 and *Erg*+149 (Fig. 1a), of which
93 only the *Erg*+85 region had previously been tested in transgenic mice (16). Inspection of the
94 gene loci of all nine TFs resulted in the identification of 35 candidate *cis*-regulatory regions
95 (Fig. 1b, Fig. 1-figure supplements 1-8). In addition to the 14 haematopoietic enhancers
96 previously published, eight of the 35 new candidate regulatory regions had previously been
97 shown not to possess activity in tissues of the blood system of mouse embryos: *Gata2*-83
98 (*Gata2*-77), *Gfilb* promoter (20), *Spi1*-18, *Spi1* promoter (22), *Runx1* P1 promoter (29),
99 *Tall*-9, *Tall* promoter (30) and *Tall*+6 (31). Of the remaining 27 candidate *cis*-regulatory
100 regions, two coincided with genomic repeat regions (*Runx1*-322 and *Runx1*+1) and were
101 excluded from further analysis because mapping of ChIP-Seq reads to such regions is
102 ambiguous. Since a comprehensive understanding of regulatory interactions among the nine
103 HSPC TFs requires *in vivo* validation of candidate regulatory regions, we next tested the
104 remaining 25 candidate *cis*-regulatory regions for their ability to mediate reporter gene
105 expression in embryonic sites of definitive haematopoietic cell emergence and colonisation,
106 namely the dorsal aorta and foetal liver of E10.5 to E11.5 transgenic *LacZ*-reporter mouse
107 embryos. For the *Erg* locus, this analysis revealed that in addition to the previously known
108 *Erg*+85 enhancer, the *Erg*+65 and *Erg*+75 regions also displayed activity in the dorsal aorta
109 and/or the foetal liver, while the *Erg*+90 and *Erg*+149 regions did not (Fig. 1c). Careful

110 inspection of a total of 188 transgenic mouse embryos revealed that nine of the 25 identified
111 regions showed *LacZ* expression in the dorsal aorta and/or foetal liver (Fig. 1b, Fig. 1-figure
112 supplements 1-9). This large scale transient transgenic screen therefore almost doubled the
113 number of known *in vivo* validated early haematopoietic regulatory elements for HSPC TFs.

114

115 **ChIP-Seq maps for a second progenitor cell line validate core regulatory interactions**

116 Although HPC7 cells are a useful model cell line for the prediction of genomic regions with
117 haematopoietic activity in transgenic mouse assays (16), they are refractory to most gene
118 transfer methods and therefore not suitable for functional characterisation of regulatory
119 elements using standard transcriptional assays. By contrast, the 416b myeloid progenitor cell
120 line can be readily transduced by electroporation and therefore represents a candidate cell line
121 for functional dissection of individual regulatory elements. As ChIP-Seq profiles in 416b cells
122 had not been reported previously, we performed ChIP-Seq for H3K27ac and the nine TFs in
123 this cell line (Fig. 2a, Fig. 2-figure supplements 1-8). Alongside with our previously published
124 HPC7 data, this new 416b dataset represents the most complete genome-scale TF-binding
125 analysis in haematopoietic progenitor cell lines to date, with all new data being freely
126 accessible under the following GEO accession number GSE69776 and also at
127 <http://codex.stemcells.cam.ac.uk/>. Genome-wide TF binding patterns in 416b and HPC7 cells
128 were closely related when compared with binding profiles for the same factors in other
129 haematopoietic lineages (Fig. 2b, Fig. 2-figure supplement 9). Inspection of the gene loci for
130 the nine HSPC TFs not only revealed many similarities between 416b and HPC7 cells, but also
131 some differences in TF binding patterns. Overall, TF occupancy at the 23 regions with activity
132 in haematopoietic tissues (14 previously published (16-26) and 9 newly identified) does not
133 change between the two cell types in 71 % of all cases (147 of 207 binding events), is gained in

134 416b cells in 16 % (33 of 207) and lost in 13 % (27 of 207) of cases compared to HPC7 cells
135 (Fig. 2c). Next, all 23 elements were filtered to only retain those elements which were bound
136 by at least 3 of the 9 TFs and displayed elevated H3K27ac in HPC7 and 416 cells. This led to
137 the removal of the *Gata2-3*, which is not bound by any of the nine TFs in either cell type,
138 *Gata2-92* and *Gfilb+13*, which are only bound by one or no TFs in 416b cells, and *Fli1-15*,
139 which is not acetylated in 416b cells (Fig. 2c, Fig. 2-figure supplements 1-3). Overall, 19 *cis*-
140 regulatory regions were therefore taken forward as a comprehensively validated set of regions
141 for the reconstruction of an HSPC regulatory network model.

142

143 **Comprehensive TFBS mutagenesis reveals enhancer-dependent effects of TF binding**

144 The reconstruction of a core regulatory network model requires information about the effect of
145 TF binding on gene expression, which can be activating, repressing or non-functional. In order
146 to analyse the effects of all TF binding events at all 19 regulatory regions, we performed
147 luciferase reporter assays in stably transfected 416b cells. Based on multiple species
148 alignments between five species (mouse, human, dog, platypus, opossum), we identified
149 conserved TFBSs for the nine TFs (Fig. 3a, Fig. 3-figure supplements 1-19), and generated
150 mutant constructs for each of the 19 regulatory regions, resulting in 87 reporter constructs that
151 were tested by luciferase assays (19 wild-type, 68 mutants). To ensure that DNA binding of the
152 TFs was abrogated, the key DNA bases involved in DNA-protein interactions were mutated
153 and the resulting sequences were scanned to ensure that no new binding sites were created (32).
154 For each of the 19 regulatory regions the conserved TFBSs were mutated by family, for
155 example all six Ets sites within the *Erg+65* region were mutated simultaneously in one
156 construct and this element was then treated as the *Erg+65_Ets* mutant. TFBS mutations
157 reduced or increased activity compared to the wild-type enhancer, or indeed had no significant

158 effect (Fig. 3b, Fig. 3-figure supplements 1-18). For instance, at the *Erg*+65 region, mutation of
159 the six Ets binding sites or the three Gata binding sites reduced luciferase activity, whereas
160 mutation of the three Ebox or the three Gfi motifs increased luciferase activity (Fig. 3b).
161 Comparison of the luciferase assay results for all 19 *cis*-regulatory regions (Fig. 3c) reveals
162 that for each motif class mutation can result in activation, repression or no-change. This
163 observation even extends to single gene loci, where for example mutation of the Gata site
164 reduced activity of the *Erg*+65 region, but increased activity of the *Erg*+85 enhancer (Fig. 3c).
165 Taken together, this comprehensive mutagenesis screen highlights the dangers associated with
166 extrapolating TF function simply from ChIP-Seq binding events and thus underlines the
167 importance of functional studies for regulatory network reconstruction.

168

169 **Dynamic Bayesian network modelling can incorporate complex regulatory information**
170 **and shows stabilization of the HSPC expression state**

171 We next set out to construct a regulatory network model that incorporates the detailed
172 regulatory information obtained for potential cross-regulation of the nine HSPC TFs obtained
173 in the previous sections (summarised in Fig. 4a). We focussed on three categories of causal
174 relationships: (i) one or several TFs can bind to a certain type of motif at a given regulatory
175 region, and the probability of a motif being bound depends on the expression levels of the
176 relevant TFs; (ii) TFBS mutations at a given regulatory region altered luciferase activities
177 compared to the wild-type, thus capturing the impact of TF binding on activity of the given
178 regulatory region; (iii) individual regulatory regions show varying degrees of activation over
179 baseline controls, which translate into different relative strengths of individual *cis*-regulatory
180 regions. To incorporate this multi-layered experimental information, we constructed a three-
181 tier dynamic Bayesian network (DBN) to jointly represent all those causal relationships (see

182 Material and Methods and Fig. 4b). The reconstructed DBN represents a first-order time-
183 homogeneous Markov process, which is a stochastic process where the transition functions are
184 the same throughout all time points and the conditional probability distribution of future states
185 depends only on the present state (see Material and Methods). The model is calculated so that
186 the expression at $t+1$ is influenced by the expression at $t0$; analogously, the expression at $t0$ is
187 influenced by the expression at $t-1$, and so on. Therefore, though the model does not
188 incorporate “epigenetic memory”, past expression levels directly influence current expression
189 levels. Model execution therefore permits the simulation of gene expression states in single
190 cells over time, as well as the calculation of gene expression distributions for each gene across
191 a population of simulated single cells.

192 Having generated a DBN model incorporating extensive experimental information, we next
193 investigated the expression states following model execution. First, we investigated whether
194 the network model was compatible with the HSPC expression profile from which all the
195 experimental data are derived, namely co-expression of all nine TFs. To this end, model
196 execution was initiated with expression levels for all nine TFs set at the midpoint level of 0.5.
197 A representative single cell modelled over time rapidly adopts characteristic levels of
198 expression for each of the nine genes, with some genes showing perpetual fluctuations (Fig.
199 4c). The same expression levels were reached when the model was initiated with expression
200 starting at 0.2, 0.8 or with initially only FLI1, RUNX1 and TAL1 being expressed at 0.5 (Fig.
201 4-figure supplement 1). We next modelled the overall distribution of the nine TFs as might be
202 seen in a cell population by running 1000 model simulations (Fig. 4d). This analysis
203 demonstrated that our model is compatible with co-expression of all nine genes within the
204 same single cell. Moreover, stable expression over time for some genes as well as oscillations
205 around a characteristic mean expression level for other genes suggests that our model may have

206 captured those aspects of HSPC regulatory networks that ensure maintenance of
207 stem/progenitor cells.

208

209 **Relative stability to experimental perturbation is recapitulated by the model**

210 The TFs TAL1 and LYL1 are important regulators of adult haematopoiesis, but the deletion of
211 each factor individually has only minor effects on adult HSC function (33-35). Combined
212 deletion in adult HSCs however causes a severe phenotype with rapid loss of HSPCs (36). We
213 wanted to investigate to what extent our computer model could recapitulate these known
214 phenotypes through *in silico* perturbation simulations. To quantify if a change in the expression
215 profile of a given TF was significant we performed a Wilcoxon rank-sum test. Interestingly,
216 this significance calculation demonstrated that both large and small fold-changes can be
217 significant. Simulated perturbation of just LYL1 caused significant alterations to the expression
218 profiles of *Gfi1b*, *Tal1*, *Fli1* and *Gata2*, but none of these were associated with a substantial
219 shift in mean expression levels (Fig. 5a, Fig. 5-figure supplement 1). Perturbation of just TAL1
220 caused significant changes to the expression profiles of *Runx1*, *Gfi1b* and *Gata2*, and again
221 none of these were associated with a substantial shift in expression levels (Fig. 5b, Fig. 5-figure
222 supplement 1). Simultaneous deletion of both factors caused significant changes in gene
223 expression profiles in all TFs except for *Fli1*. Unlike for the single TF perturbations, *Gata2* and
224 *Runx1* showed substantial shifts in expression levels when both LYL1 and TAL1 were
225 simulated to be knocked down (Fig. 5c, Fig. 5-figure supplement 1). Of note, the significance
226 calculations highlight that there may be no one perfect way to visualize these small fold-change
227 alterations. We therefore also generated histogram plots as an alternative visualization (Fig. 5-
228 figure supplement 2).

229 We next wanted to compare model predictions with actual experimental data in the 416b cell
230 line, from which the information for model construction had been derived. Because our DBN
231 model is particularly suited to model the expression states in single cells, we compared
232 predicted and experimentally observed effects of knockdown or overexpression in single cells.
233 To this end we knocked down the expression of TAL1 in 416b cells by transfecting the cells
234 with siRNA against *Tal1* (siTal1) or control siRNA (siCtrl). Forty-eight hours after
235 transfection, gene expression for the nine network genes was analysed in 44 siTal1 treated cells
236 and 41 siCtrl treated cells. Importantly, 29 of 44 cells (66%) transfected with siTal1 showed no
237 expression of *Tal1* anymore, demonstrating the successful knockdown (Fig. 5d, Fig. 5-source
238 data). Down-regulation of TAL1 caused a significant change in the expression profiles of
239 *Tal1*, *Fli1* and *Gfi1b*, but a substantial shift of median expression was only observed for *Tal1*
240 (Fig. 5-figure supplement 1). Experimental validation therefore confirmed the occurrence of
241 statistically significant small-fold changes in expression profiles following single TF
242 knockdown, although there was no perfect match between the genes affected in the model and
243 experiment. To extend comparisons between model predictions and experimental validation,
244 we investigated the consequences of knocking down the expression of PU.1 and
245 overexpressing GFI1B. Complete removal of PU.1 *in silico* after the model had reached its
246 initial steady state had no effect on the expression levels of the other TFs (Fig. 6a). To
247 investigate whether the model prediction is comparable to experimental data obtained from
248 single cells, single cell gene expression analysis using the Fluidigm Biomark HD platform was
249 performed using 416b cells transduced with shRNA against PU.1 (shPU.1) or luciferase
250 (shluc). Three days after transduction, 121 shPU.1 and 123 shluc transduced single cells were
251 analysed for their expression of *Spi1* and the other eight TFs of the network. 18 shPU.1-
252 transduced cells (15%) showed a complete loss of *Spi1*, and expression of *Spi1* in the
253 remaining cells was markedly reduced compared to the control cells (shluc) (Fig. 6a, Fig. 5-

254 source data), highlighting the efficiency of the PU.1 knockdown. *Spi1*, *Runx1*, *Erg* and *Fli1*
255 showed a significant change in expression profiles after the depletion of PU.1, but this involved
256 a substantial shift in median expression levels only for *Spi1* and *Runx1* (Fig. 5-figure
257 supplement 1). Expression profiles of the remaining five TFs did not change as a result of
258 reduced PU.1 levels (Fig. 6a, Fig. 5-source data), therefore mostly confirming the model
259 prediction.

260 Next, we modelled GFI1B overexpression *in silico* by increasing the expression level of *Gfi1b*
261 to the maximum value after the model had reached its initial steady state which led to a
262 significant change in the expression profiles of *Gfi1b*, *Meis1*, *Erg* and *Runx1*, although only
263 *Gfi1b* and *Meis1* showed a substantial shift in median expression levels (Fig. 6b, Fig. 5-figure
264 supplement 1, Fig. 5-source data). Expression profiles of the other five TFs were unaltered.
265 Single cell gene expression analysis of 90 single 416b cells transduced with a *Gfi1b*-expressing
266 vector and 104 single 416b cells transduced with an empty control vector showed a significant
267 increase in the expression of *Gfi1b* and a significant alteration to the expression profile of *Erg*,
268 but only the changes to *Gfi1b* involved a substantial shift in median expression levels. No
269 significant expression changes were seen in any of the other seven network genes (Fig. 6b).
270 Both PU.1 and GFI1B perturbation studies therefore emphasize the resilience of the HSPC TF
271 network to single TF perturbation. Moreover, our *in silico* model reflects this, thus suggesting
272 that the comprehensive experimental information used to construct the network model has
273 allowed us to capture key mechanistic aspects of HSPC regulation. Of note, there were no
274 short-term major expression changes immediately after the perturbation in the *in silico*
275 simulations for the three single TF perturbation described above. For completeness we
276 performed *in silico* modelling for all permutations of single TF knockdown / overexpression as
277 well as all pairwise combinations of all 9 TFs analysed (a total of 162 simulations, Fig. 6-figure
278 supplement 1).

279

280 **Major perturbations by the AML-ETO oncoprotein are captured by the network model**

281 As the TF network described above is relatively stable to single TF perturbations, we set out to
282 test whether a simulation that mimics the situation present in leukemic cells can influence the
283 expression states of the nine TFs in our network. The *Aml-Eto9a* translocation is amongst the
284 most frequent mutations in AML (reviewed in (37)). The resulting fusion protein is thought to
285 act in a dominant negative manner to repress RUNX1 target genes. To simulate the leukemic
286 scenario caused by AML-ETO expression, we fixed the level of *Runx1* to be the maximum
287 value 1 and at the same time converted all activating inputs of RUNX1 to inhibiting inputs in
288 our DBN model. Interestingly, this simulation of a “leukemic” perturbation caused significant
289 expression changes to all eight of the core HSPC TFs (Fig. 6c). To compare the AML-ETO
290 simulation results with experimental data, we utilised a doxycycline-inducible expression
291 system to generate 416b cells with inducible expression of AML1-ETO fused to a mCherry
292 reporter via a self-cleaving 2A peptide spacer. Following doxycycline induction, 56 single
293 mCherry positive and 122 single mCherry negative 416b cells were analysed by single cell
294 gene expression. Significant gene expression changes can be seen in six of the nine core HSPC
295 TFs (all except *Tal1*, *Erg* and *Gata2*) thus highlighting significant overlap between predictions
296 and experimental validation, although there are also notable differences between model
297 predictions and the experimental data (see for example *Gata2*; Fig. 6c, Fig. 5-figure supplement
298 1, Fig. 5-source data). These results demonstrate that our new HSPC network model can
299 capture many gene expression changes caused by ectopic expression of a leukemia oncogene as
300 well as providing a useful model for normal HSPC transcriptional regulation. The inability of
301 any model to completely recapitulate experimental data is not unexpected. Possible reasons in
302 our case may include more complex activities of the onco-fusion protein than would be
303 captured by our assumption that its “only” function is as a straightforward dominant-negative

304 effect, or the fact that the computational model is a closed system of only the 9 network TFs,
305 whereas the experimental single cell perturbation is subject to possible knock-on consequences
306 from gene changes outside of the 9-TF network.

307

308 **Discussion**

309 Transcription factor networks are widely recognised as key determinants of cell type identity.
310 Since the functionality of such regulatory networks is ultimately encoded in the genome, the
311 logic that governs interactions between network components should be identifiable, and in due
312 course allow for the construction of network models that are capable of capturing the behaviour
313 of complex biological processes. However, the construction of such network models has so far
314 been severely restricted because the identification and subsequent functional characterisation of
315 mammalian regulatory sequences represent major challenges, and the connectivity and
316 interaction rules within regulatory networks can be highly complex. Here we report a
317 comprehensive mammalian transcriptional network model that is fully grounded in
318 experimental data. Model simulation coupled with subsequent experimental validation using
319 sophisticated single cell transcriptional assays revealed the mechanistic basis for cell state
320 stability within a haematopoietic progenitor model cell line, and also how a leukemogenic TF
321 fusion protein can perturb the expression of a subset of key blood stem cell regulators.

322 Pictorial representations of putative network models are commonly shown in publications
323 reporting ChIP-Seq TF binding datasets (38). However, due to the lack of experimental
324 underpinning, such representations are simple images that do not encode any of the underlying
325 gene regulatory logic, and importantly therefore cannot provide executable computational
326 models that can be used to simulate biological systems. Although the experimentally-grounded
327 network model shows good agreement with the relative expression states of the nine TFs for

328 the wild-type as well as the perturbation data, model predictions are not correct in all cases.
329 Apart from the obvious caveat that any computer model is an abstraction of reality and
330 therefore will not be correct in every detail, it also needs to be stressed that we treat the nine
331 TFs as an isolated module for the computer simulations, and therefore could not account for
332 possible influences by additional genes that may affect single cell gene expression
333 measurements in the perturbation experiments.

334 Statistical significance calculations demonstrated that both the computer model and the
335 experimental data showed significant changes in gene expression profiles that were associated
336 with minimal fold-change alterations to median expression levels. Such alterations to
337 expression profiles were prevalent in both single and double-gene perturbations, whereas
338 substantial shifts in median expression were mostly restricted to the double perturbations (and
339 also the AML-ETO oncogene overexpression). This observation suggests that (i) our approach
340 has the capacity to reveal aspects of the fine-grained nature of biological networks, and (ii) the
341 network presented in this study is largely resistant to perturbations of individual TFs in terms
342 of substantial fold-change alterations in median expression levels. We believe that it may well
343 be possible that the statistically significant small-fold changes in HSPC network genes may be
344 responsible for the mild phenotypes seen when major HSPC regulators are deleted in adult
345 HSPCs. *Tal1*^{-/-} mice for example are not viable because TAL1 is absolutely required for
346 embryonic blood development (39), yet deletion of TAL1 in adult HSCs only causes minor
347 phenotypes (33). Another noteworthy observation is that it would have been impossible to
348 detect the statistically significant yet small fold-changes using conventional expression
349 profiling, because they only become apparent following the statistical analysis of expression
350 distributions generated by assaying lots of single cells. More generally it is important to
351 acknowledge that the question of how close the present model comes to capturing the

352 underlying biological processes can only be revealed through much more exhaustive
353 experimental validation studies.

354 A potential caveat for network reconstruction based on identification of regulatory elements
355 comes from the difficulties associated with capturing negative regulatory elements. As shown
356 elegantly for CD4 and CD8 gene silencers in the lymphoid lineage, TFs involved in the early
357 repression of a locus are not required for maintenance of the silenced state (40, 41).
358 Identification of negative regulatory inputs may therefore require an expansion of datasets to
359 look across sequential developmental stages. It will therefore be important in the future to
360 extend the work presented here to include additional HSPC regulators as well as additional
361 stages along the haematopoietic differentiation hierarchy. Of note, TF-mediated cellular
362 programming experiments have demonstrated that modules of 3-4 TFs are able to confer cell-
363 type specific transcriptional programmes (42-45), consistent with the notion that a network
364 composed of nine key HSPC regulators is able to capture useful information about HSPC
365 regulatory programmes.

366 One of the most striking observations of the regulatory network defined here is the high degree
367 to which the HSPC expression state is stabilised. As such, this model is different from
368 previous experimentally-grounded transcriptional regulatory network models (46). These
369 earlier model organism networks have inherent forward momentum, where the model captures
370 the progression through successive embryonic developmental stages characterised by distinct
371 expression states.

372 The model reported here is based on and validated with data from haematopoietic progenitor
373 cell lines, which can differentiate (47, 48), but can also be maintained in stable self-renewing
374 conditions. A recent study by Busch and colleagues tracked labelled Tie2⁺ HSCs in the bone
375 marrow, and showed that haematopoietic progenitors *in vivo* are also characterised by a

376 substantial self-renewal capability, therefore highlighting the stable state in which they can
377 reside for several months (49). The observed stability of the HSPC expression state presented
378 here is therefore likely to capture aspects of the regulatory mechanisms maintaining the steady
379 state of primary haematopoietic progenitor cells, a notion reinforced further by the fact that our
380 model is based on *in vivo* validated regulatory elements.

381 The two types of models therefore accurately capture the properties of the distinct biological
382 processes, e.g. driving developmental progression on the one hand, and maintaining a given
383 cellular state on the other. Different design principles are likely to be at play, with feed-
384 forward loops representing key building blocks of early developmental GRNs, while the
385 network described here shows an abundance of auto-regulatory feedback loops and partially
386 redundant enhancer elements, both of which may serve to stabilise a given cellular state.

387 Of particular interest may be the organisation of the *Runx1* gene locus, where RUNX1 protein
388 provides positive feedback at some, and negative feedback at other HSPC enhancers. Given
389 that these different enhancers employ overlapping yet distinct sets of upstream regulators, it is
390 tempting to speculate that such an arrangement not only stabilises a given expression level, but
391 also provides the means to either up- or down-regulate RUNX1 expression in response to
392 diverse external stimuli that may act on specific RUNX1 co-factors at either the repressing or
393 activating RUNX1 binding events. Taken together, we report widely applicable experimental
394 and computational strategies for generating fully validated regulatory network models in
395 complex mammalian systems. We furthermore demonstrate how such a model derived for
396 blood stem/progenitor cells reveals mechanisms for stabilisation of the progenitor cell state,
397 and can be utilised to analyse core network perturbations caused by leukemic oncogenes.

398

399 **Materials and Methods**

400 **ChIP-Sequencing and data processing**

401 The mouse myeloid progenitor 416b cell line (48) was received from Chester Beatty lab and
402 confirmed to be mycoplasma free. The cells were cultured in RPMI with 10 % FCS and 1 %
403 Penicillin/Streptomycin.

404 ChIP assays were performed as previously described (16, 27), amplified using the Illumina
405 TruSeq ChIP Sample Prep Kit and sequenced using the Illumina HiSeq 2500 System following
406 the manufacturer's instructions. Sequencing reads were mapped to the mm10 mouse reference
407 genome using Bowtie2 (50), converted to a density plot and displayed as UCSC genome
408 browser custom tracks. Peaks were called using MACS2 software (51). Mapped reads were
409 converted to density plots and displayed as UCSC genome browser custom tracks. The raw and
410 processed ChIP-Seq data have been submitted to the NCBI Gene Expression Omnibus
411 (www.ncbi.nlm.nih.gov/geo) and assigned the identifier GSE69776. A binary binding matrix
412 was created using in-house scripts, clustered using the dice coefficient and a heatmap was
413 plotted using gplots in R in order to compare newly generated ChIP-Seq data with previously
414 published data (52).

415

416 **Analysis of enhancer activity in transient transgenic mouse embryos**

417 Genomic fragments spanning the candidate *cis*-regulatory regions were generated by PCR or
418 ordered as gBlocks (Life Technologies GmbH) and cloned downstream of the LacZ gene in an
419 hsp68LacZ (Runx1 constructs) or SVLacZ (all other constructs) reporter vector. Coordinates of
420 candidate chromosomal regions and corresponding primer sequences are given in Fig. 3-figure
421 supplement 20. For Runx1, E10 mouse transient transgenic embryos carrying LacZ enhancer-
422 reporter constructs were generated by pronuclear injection of (C57BL/6 x CBA)/F2 zygotes
423 following standard procedures. Transgenic embryos were identified by LacZ-specific PCR on

424 genomic DNA isolated from yolk sac (5'-GCAGATGCACGGTTACGATG-3'; 5'-
425 GTGGCAACATGGAAATCGCTG-3'). Xgal staining and cryostat sectioning were performed
426 as previously described (23). Embryos were photographed using a Leica MZFLIII microscope,
427 Leica DFC 300F digital camera (Leica Microsystems, Milton Keynes, UK) and Openlab
428 software (Improvision, Coventry, UK) and sections were examined using a Nikon Eclipse
429 E600 microscope (Nikon, Japan) equipped with 20x and 40x Nomarski objectives.
430 Photographs were taken using a Nikon DXM 1200c Digital Camera (Nikon, Tokyo, Japan).
431 E11.5 transient transgenic embryos of all other candidate *cis*-regulatory regions were generated
432 by Cyagen Biosciences Inc (Guangzhou, China). Whole-mount embryos were stained with 5-
433 bromo-4-chloro-3-indolyl- β -d-galactopyranoside (X-Gal) for β -galactosidase expression and
434 photographed using a Nikon Digital Sight DS-FL1 camera attached to a Nikon SM7800
435 microscope (Nikon, Kingston-upon-Thames, UK). Candidate transgenic mouse embryos with
436 LacZ staining in haematopoietic tissues were subsequently embedded in paraffin, stained with
437 0.1 % (w/v) Neutral Red and cut into 6 μ m deep longitudinal sections. Images of sections were
438 acquired with a Pixera Penguin 600CL camera attached to an Olympus BX51 microscope. All
439 images were processed using Adobe Photoshop (Adobe systems Europe, Uxbridge, United
440 Kingdom).

441

442 **Luciferase reporter assays**

443 Wild-type and mutant DNA fragments for candidate regulatory regions were either cloned
444 using standard recombinant DNA techniques, ordered as gBlocks (Life Technologies) or
445 obtained from GeneArt® by Life Technolgies. DNA fragments were cloned into pGL2 basic or
446 pGL2 promoter vectors from Promega using restriction enzymes or by Gibson Assembly.
447 TFBSs for the nine TFs of interest (corresponding DNA sequences are listed in Fig. 3-figure
448 supplement 19) were identified based on multiple species alignments between five species

449 (mouse, human, dog, platypus, opossum). Where a region contained multiple instances of the
450 same motif, a single mutant construct with all relevant motifs mutated simultaneously was
451 generated (for generated point mutations check Fig. 3a and Fig. 3-figure supplements 1-18).
452 Where TF binding was observed in ChIP-Seq experiments in 416b cells, but the TFBS was not
453 conserved, the motifs present in the mouse sequence were mutated. Stable transfections of the
454 416b cell line were performed using 10 µg reporter construct, 2 µg neomycin resistance
455 plasmid and 1×10^7 416b cells in 180 µl culture medium per pulse. The sample was
456 electroporated at 220V and capacitance of 900 µF using the GenePulser Xcell Electroporation
457 System (Bio-Rad). Immediately after transfection, the sample was split into four culture plates.
458 Twenty-four hours after transfection Geneticin G418 (Gibco by Life Technologies) at a final
459 concentration of 0.75 mg/ml was applied to the culture to select for transfected cells. The
460 activity of the luciferase reporter constructs was measured 12-16 days after transfection by
461 using a FLUOstar OPTIMA luminometer (BMG LABTECH). The luciferase activity was
462 normalised to the cell number and presented as relative activity compared to the wild-type
463 construct. All assays were performed at least three times in quadruplicates.

464

465 **Single cell gene expression and data analysis**

466 The TAL1 knockdown was performed using pools of siRNA against Tal1 (Dharmacon) which
467 were transfected into 416b cells. Briefly, 1×10^6 cells were electroporated with either a control
468 or Tal1 siRNA. Forty-eight hours after transfection, cells were sorted into 96 well PCR plates
469 containing lysis buffer using the BD Influx Cell Sorter.

470 The PU.1 knockdown was performed as previously described (27).

471 The MigR1-Gfi1b retroviral expression vector and the corresponding empty vector control (53)
472 were used for GFI1B overexpression. Two million 416b cells were transduced with the above

473 listed vectors by adding viral supernatant and 4 µg/ml polybrene to the cells, followed by
474 centrifugation at 900 x g for 90 min at 32°C and incubation with 5% CO₂ at 32°C. Half of the
475 media was then replaced with fresh culture media and cells were incubated at 37°C with 5 %
476 CO₂. Forty-eight hours after transduction, GFP⁺ cells for each cell population were sorted into
477 96 well PCR plates containing lysis buffer using the BD Influx Cell Sorter.

478 To induce AML1-ETO9a expression, the 416b cell line was co-transfected with: 1) a plasmid
479 containing the tetracycline transcription silencer (tTS), the tetracycline transactivator (rtTA)
480 and blasticidine resistance under the control of the *EF1α* promoter; 2) a plasmid containing the
481 entire *Aml-Eto9a* cDNA (obtained from vector MigR1-AE9a, Addgene no. 12433) in frame
482 with a F2A element and the mCherry protein under the control of a tetracycline responsive
483 element; and 3) transposase PL623 (54) (kindly donated by Pentao Liu, Sanger Institute,
484 Cambridge) in order to promote simultaneous stable integration of the two constructs described
485 above. After 6 days of culture without selection, cells were incubated with 1 µg/ml of
486 Doxycycline for 24 hours and then stained with DAPI. mCherry positive and negative cells that
487 did not stain with DAPI were sorted into 96 well PCR plates containing lysis buffer using the
488 BD Influx Cell Sorter.

489 Single cell gene expression analysis was performed using the Fluidigm BioMark platform
490 followed by bioinformatics analysis as previously described (20). All cells that express less
491 than 48 % of genes assayed were removed from the analysis for PU.1 knockdown and GF11B
492 over-expression, all cells expressing less than 56 % of genes assayed were removed from the
493 TAL1 knockdown and all cells that express less than 44 % of genes assayed were removed
494 from the analysis for the AML-ETO9a induction. Importantly, this thresholding resulted in the
495 removal of similar numbers of cells in both the perturbation and control arms of the
496 experiments. The raw data as well as the normalised data (normalised to Ubc and Polr2a) of the
497 gene expression analysis are listed in Fig. 5-source data).

498

499 **Computational modelling**

500 The first-order DBN shown in Fig. 4b was established on the basis of regulatory information
501 summarized in Fig. 4a. The DBN essentially presents a discrete-time stochastic process that
502 has the Markov property, i.e. the state of the process at the next time point depends purely on
503 its state at the current time point. Also note that this is a time-homogeneous (or time-invariant)
504 DBN, where the transition functions/matrices are the same throughout all time points.

505 To specify parameters of the DBN, we defined a motif family at a specific regulatory region as
506 a unique binary variable; with value “1” indicating that no motif of a motif family is bound at
507 the specific region and value “2” indicating that at least one motif of the motif family is bound
508 by a TF at this region. We assumed that any of the following three factors can lead to a higher
509 probability of a motif being bound by a TF and therefore taking the value 2: (i) more motifs of
510 the same type present within a regulatory region; (ii) multiple TFs that can bind to the same
511 motif, such as TAL1 and LYL1 both binding to Eboxes; (iii) higher expression levels of the
512 TFs. The probabilities were thus calculated based on these three sources of information (see
513 below for an example). We next defined that every regulatory region was a continuous variable
514 on the close interval $[0, 1]$, and its value was determined by the accumulated effects of all
515 motifs present with the regulatory region. Finally, the expression levels of the nine TFs were
516 also defined as continuous variables ranging from 0 to 0.8, and their expression levels were
517 determined by the accumulated activities of the relevant regulatory regions.

518 Considering that variables in the top tier of the DBN are binary whereas those in the middle
519 and bottom tiers are continuous, we found conditional linear Gaussian distribution (55) to be an
520 appropriate generic representation of the intra-slice conditional probability distributions.
521 Specifically, the regression coefficient of a regulatory region on a motif family was estimated
522 by normalizing the logarithmic deviation of luciferase activity, where deviation refers to the

523 change of luciferase activity between the wild-type and the mutated (one motif family at a
524 time) regulatory region (see below for a demonstration). Using the logarithmic deviation
525 allowed us to account for the differences in effect sizes of various motif families by rescaling
526 the differences to a comparable range. Similarly, for each of the nine genes, the regression
527 coefficient of its expression level on a relevant regulatory region was estimated by normalizing
528 the logarithmic deviation of luciferase activity, where deviation refers to the change of
529 luciferase activity compared to the empty vector controls. All Matlab source codes are
530 available at <https://github.com/Huange> and also <http://burrn-sim.stemcells.cam.ac.uk/>.

531

532 Detailed explanation of the modelling of each tier of the DBN:

533 *a) Estimating the discrete probability distribution of a motif variable*

534 The probability of a motif family at a given regulatory region taking value 1 or 2 (i.e., being
535 unbound or bound) was calculated based on: (i) the number of such motifs in that regulatory
536 region; (ii) the expression levels of the relevant TFs.

537 For example, three Ebox motifs were found at *Erg+65* (Fig. 3a). They can be bound by either
538 TAL1 or LYL1. Thus, we assigned that $P(\text{Ebox}@Erg+65=1)$ and $P(\text{Ebox}@Erg+65=2)$ were
539 determined by $\{3, \text{TAL1}, \text{LYL1}\}$. We assumed that (i) the expression level of a TF is
540 proportional to the probability of that TF binding to a target motif; and (ii) the bindings of TFs
541 to multiple motifs are independent events. Gene expression levels were defined within the
542 closed interval $[0, 1]$, which is identical to the possible range of probabilities. For ease of
543 calculation, we took the expression level of a TF as its probability of binding to a motif.

544 Accordingly, we have

$$545 \quad \tilde{P}(\text{Ebox}@Erg+65=1) = (1-p)^3 \times (1-q)^3 \quad (1)$$

$$546 \quad \tilde{P}(\text{Ebox}@Erg+65=2) = \sum_{n=1}^3 C(3, n) \times p^n \times (1-p)^{(3-n)} \times (1-q)^3$$

547
$$+ \sum_{n=1}^3 C(3, n) \times q^n \times (1-q)^{(3-n)} \times (1-p)^3 \quad (2)$$

548
$$+ \sum_{n=1}^2 \sum_{m=1}^{3-n} C(3, n) \times p^n \times (1-p)^{(3-n)} \times C((3-n), m) \times q^m \times (1-q)^{(3-m)}$$

549 where p and q represent the expression levels of TAL1 and LYL1, respectively.

550 However, to remove the bias introduced by simply taking the expression level of a TF as its
551 probability of binding to a motif, we further normalized the resulting probabilities as below:

552
$$\tilde{Z} = \tilde{P}(\text{Ebox@Erg+65} = 1) + \tilde{P}(\text{Ebox@Erg+65} = 2) \quad (3)$$

553
$$P(\text{Ebox@Erg+65} = 1) = \tilde{P}(\text{Ebox@Erg+65} = 1) / \tilde{Z} \quad (4)$$

554
$$P(\text{Ebox@Erg+65} = 2) = \tilde{P}(\text{Ebox@Erg+65} = 2) / \tilde{Z} \quad (5)$$

555 It should be mentioned that the number of the same motifs in a regulatory region was directly
556 taken into account in the estimation of probabilities. One may raise the question of whether this
557 number has such strong power. Specifically, should the exponents in equations (1) and (2)
558 change linearly, or less than linearly, along with the increase in the number of Ebox motifs? To
559 address this issue, we replaced all exponents with their square roots and rerun the whole set of
560 simulations (data not shown). Results showed that using the square roots instead of the original
561 numbers (i) caused a more evenly distributed expression of the nine TFs over the hypothetical
562 interval [0, 1], (ii) captured the same trend in gene expression changes in some perturbations
563 (e.g. the AML-ETO simulation), but (iii) led to decreased expression levels of certain TFs in
564 other perturbations (e.g. PU.1 knockdown and GFI1B over-expression), which therefore
565 disagrees with the experimental data. In order to capture a better agreement of computational
566 and experimental results, we directly used the number of motifs to estimate the discrete
567 probability distributions.

568

569 *b) Estimating the activity of a regulatory region*

570 The regression coefficient of a regulatory region on a motif family was estimated by
571 normalizing the logarithmic deviation of luciferase activity, e.g. comparing the change of
572 luciferase activity between the wild-type and mutated constructs. For example, when the
573 luciferase activity for the wild-type *Erg+65* region was set to 100 %, the simultaneous
574 mutation of all Ebox or Gfi motifs at this region resulted in increased luciferase activity (181.2
575 % or 475.9 %, respectively) (Fig. 3b). In contrast, simultaneous mutation of all Ets or Gata
576 motifs at this region led to reduced luciferase activity (1.3 % or 14.5 %, respectively). Based on
577 this information, we estimated the regression coefficient of the *Erg+65* region on a relevant
578 motif family in the following way:

$$579 \quad \alpha_i = \log\left(\frac{100}{l_k}\right) \times \left(\sum_k \left|\log\left(\frac{100}{l_k}\right)\right|\right)^{-1} \quad (6)$$

580 where $k \in \{1, \dots, 4\}$, $l_1 = 181.2$, $l_2 = 475.9$, $l_3 = 1.3$, $l_4 = 14.5$; accordingly, $\alpha_1 = -0.070$,
581 $\alpha_2 = -0.185$, $\alpha_3 = 0.515$, $\alpha_4 = 0.230$. We can then formulate a linear regression equation as
582 below:

$$583 \quad \tilde{y} = \alpha_1 x_1 + \alpha_2 x_2 + \alpha_3 x_3 + \alpha_4 x_4 \quad (7)$$

584 where \tilde{y} denote the estimated luciferase activity of *Erg+65*, and x_1 , x_2 , x_3 and x_4 represent
585 the binding status of Ebox, Gfi, Ets and Gata motifs at *Erg+65*.

586 However, the minimum and maximum \tilde{y} obtained by the above formula are 0.235 (when
587 $x_1 = x_2 = 2$ and $x_3 = x_4 = 1$) and 1.235 (when $x_1 = x_2 = 1$ and $x_3 = x_4 = 2$). To make the values
588 of \tilde{y} fall in the desired closed interval $[0, 1]$, an intercept of -0.235 has to be introduced into the
589 linear regression model. In addition, a disturbance term has been included in the model in order
590 to satisfy the generic assumption of conditional linear Gaussian distribution. Finally, the fully
591 defined linear regression model regarding *Erg+65* is given as:

$$592 \quad \tilde{y} = c + \alpha_1 x_1 + \alpha_2 x_2 + \alpha_3 x_3 + \alpha_4 x_4 + \varepsilon \quad (8)$$

593 where $c = -0.235$, $\varepsilon \sim N(0, \sigma^2)$, and σ should be a very small value.

594

595 *c) Estimating the expression level of a gene*

596 For each gene studied, the regression coefficient of its expression level on a relevant regulatory
597 region was estimated by normalizing the logarithmic deviation of luciferase activity, where
598 deviation refers to the change of luciferase activity compared to an empty vector control.

599 For example, when setting the luciferase activity of the wild-type constructs to 100 %, the
600 luciferase activity of the empty vector controls relative to *Erg+65*, *Erg+75* and *Erg+85* wild-
601 types are 1.9 %, 1.0 % and 15.2 %, respectively (Fig. 3b, Fig. 3-figure supplement 1 and 2).
602 Based on these data, we estimated the expression level of *Erg* on a relevant regulatory region in
603 the following way:

$$604 \quad \beta_i = \log\left(\frac{100}{l_k}\right) \times \left(\sum_k \left|\log\left(\frac{100}{l_k}\right)\right|\right)^{-1} \quad (9)$$

605 where $k \in \{1, 2, 3\}$, $l_1 = 1.9$, $l_2 = 1.0$, $l_3 = 15.2$; accordingly, $\beta_1 = 0.379$, $\beta_2 = 0.441$,
606 $\beta_3 = 0.180$. We can then formulate a linear regression equation as below:

$$607 \quad \tilde{z} = \beta_1 \tilde{y}_1 + \beta_2 \tilde{y}_2 + \beta_3 \tilde{y}_3 \quad (10)$$

608 where \tilde{z} denote the estimated expression level of *Erg*; and \tilde{y}_1 , \tilde{y}_2 and \tilde{y}_3 represent the
609 estimated activities of *Erg+65*, *Erg+75* and *Erg+85*. Again, a disturbance term has been
610 introduced to the model in order to meet the generic assumption of conditional linear Gaussian
611 distribution. Thus, the fully defined linear regression model regarding *Erg* is given as:

$$612 \quad \tilde{z} = \beta_1 \tilde{y}_1 + \beta_2 \tilde{y}_2 + \beta_3 \tilde{y}_3 + \varepsilon \quad (11)$$

613 where $\varepsilon \sim N(0, \sigma^2)$ and σ should be a very small value.

614

615 **Statistics**

616 Significance for the results of the luciferase reporter assays was calculated by combining the p-
617 values of each experiment (generated by using the t-test function in Excel) using the Fisher's
618 method, followed by the calculation of Stouffer's z trend if necessary. Significance tests for
619 changes in TF expression levels caused by TF perturbations (both computational and
620 experimental) were evaluated by Wilcoxon rank-sum tests.

621

622 **Acknowledgments**

623 We thank the CIMR Flow Cytometry Core facility, especially Dr Chiara Cossetti, for their
624 expertise with cell sorting, Dr Marina Evangelou for her help with statistical analysis of the
625 luciferase assay data, Lucas Greder for advice on cell transfection and stimulating discussions
626 and Cyagen Biosciences and the MRC MHU Transgenic Core for generating transgenic
627 embryos. Thanks are also extended to past and present members of the Göttgens and de Bruijn
628 lab for practical assistance in the analysis of transient transgenic embryos, and to Yoram
629 Groner and Ditsa Levanon for insightful discussions. We are grateful to Barbara L. Kee
630 (University of Chicago, USA) for providing the MigRI-Gfi1b vector and Peter Laslo
631 (University of Leeds, UK) for the shPU.1 construct.

632

633 **Funding**

634 Research in the authors' laboratories was supported by Bloodwise, The Wellcome Trust,
635 Cancer Research UK, the Biotechnology and Biological Sciences Research Council, the
636 National Institute of Health Research, the Medical Research Council, the MRC Molecular
637 Haematology Unit (Oxford) core award, a Weizmann-UK "Making Connections" grant

638 (Oxford) and core support grants by the Wellcome Trust to the Cambridge Institute for Medical
639 Research (100140) and Wellcome Trust–MRC Cambridge Stem Cell Institute (097922).

640

641 **Competing Interests**

642 The authors declare that no competing interests exist.

643

644 **References**

- 645 1. Davidson EH. Emerging properties of animal gene regulatory networks. *Nature*.
646 2010;468(7326):911-20.
- 647 2. Davidson EH. Network design principles from the sea urchin embryo. *Current Opinion in*
648 *Genetics & Development*. 2009;19(6):535-40.
- 649 3. Petricka JJ, Benfey PN. Reconstructing regulatory network transitions. *Trends Cell Biol*.
650 2011;21(8):442-51.
- 651 4. Pimanda JE, Gottgens B. Gene regulatory networks governing haematopoietic stem cell
652 development and identity. *Int J Dev Biol*. 2010;54(6-7):1201-11.
- 653 5. Gottgens B. Regulatory network control of blood stem cells. *Blood*. 2015;125(17):2614-20.
- 654 6. Schutte J, Moignard V, Gottgens B. Establishing the stem cell state: insights from regulatory
655 network analysis of blood stem cell development. *Wiley interdisciplinary reviews Systems biology and*
656 *medicine*. 2012;4(3):285-95.
- 657 7. Bonzanni N, Garg A, Feenstra KA, Schutte J, Kinston S, Miranda-Saavedra D, et al. Hard-wired
658 heterogeneity in blood stem cells revealed using a dynamic regulatory network model. *Bioinformatics*
659 (Oxford, England). 2013;29(13):i80-8.
- 660 8. Krumsiek J, Marr C, Schroeder T, Theis FJ. Hierarchical differentiation of myeloid progenitors is
661 encoded in the transcription factor network. *PloS one*. 2011;6(8):e22649.
- 662 9. Narula J, Williams CJ, Tiwari A, Marks-Bluth J, Pimanda JE, Igoshin OA. Mathematical model of
663 a gene regulatory network reconciles effects of genetic perturbations on hematopoietic stem cell
664 emergence. *Developmental Biology*. 2013;379(2):258-69.
- 665 10. Narula J, Smith AM, Gottgens B, Igoshin OA. Modeling reveals bistability and low-pass filtering
666 in the network module determining blood stem cell fate. *PLoS computational biology*.
667 2010;6(5):e1000771.
- 668 11. Gazit R, Garrison BS, Rao TN, Shay T, Costello J, Ericson J, et al. Transcriptome analysis
669 identifies regulators of hematopoietic stem and progenitor cells. *Stem cell reports*. 2013;1(3):266-80.
- 670 12. Jojic V, Shay T, Sylvia K, Zuk O, Sun X, Kang J, et al. Identification of transcriptional regulators in
671 the mouse immune system. *Nature immunology*. 2013;14(6):633-43.
- 672 13. Dunn SJ, Martello G, Yordanov B, Emmott S, Smith AG. Defining an essential transcription
673 factor program for naïve pluripotency. *Science (New York, NY)*. 2014;344(6188):1156-60.
- 674 14. Zhou Q, Chipperfield H, Melton DA, Wong WH. A gene regulatory network in mouse
675 embryonic stem cells. *Proceedings of the National Academy of Sciences of the United States of*
676 *America*. 2007;104(42):16438-43.

- 677 15. Wilson NK, Foster SD, Wang X, Knezevic K, Schutte J, Kaimakis P, et al. Combinatorial
678 transcriptional control in blood stem/progenitor cells: genome-wide analysis of ten major
679 transcriptional regulators. *Cell stem cell*. 2010;7(4):532-44.
- 680 16. Wilson NK, Miranda-Saavedra D, Kinston S, Bonadies N, Foster SD, Calero-Nieto F, et al. The
681 transcriptional program controlled by the stem cell leukemia gene *Scl/Tal1* during early embryonic
682 hematopoietic development. *Blood*. 2009;113(22):5456-65.
- 683 17. Beck D, Thoms JA, Perera D, Schutte J, Unnikrishnan A, Knezevic K, et al. Genome-wide
684 analysis of transcriptional regulators in human HSPCs reveals a densely interconnected network of
685 coding and noncoding genes. *Blood*. 2013;122(14):e12-22.
- 686 18. Pimanda JE, Ottersbach K, Knezevic K, Kinston S, Chan WY, Wilson NK, et al. *Gata2*, *Fli1*, and
687 *Scl* form a recursively wired gene-regulatory circuit during early hematopoietic development.
688 *Proceedings of the National Academy of Sciences of the United States of America*.
689 2007;104(45):17692-7.
- 690 19. Wozniak RJ, Boyer ME, Grass JA, Lee Y, Bresnick EH. Context-dependent GATA factor function:
691 combinatorial requirements for transcriptional control in hematopoietic and endothelial cells. *The*
692 *Journal of biological chemistry*. 2007;282(19):14665-74.
- 693 20. Moignard V, Macaulay IC, Swiers G, Buettner F, Schutte J, Calero-Nieto FJ, et al.
694 Characterization of transcriptional networks in blood stem and progenitor cells using high-throughput
695 single-cell gene expression analysis. *Nature cell biology*. 2013;15(4):363-72.
- 696 21. Chan WY, Follows GA, Lacaud G, Pimanda JE, Landry JR, Kinston S, et al. The paralogous
697 hematopoietic regulators *Lyl1* and *Scl* are coregulated by *Ets* and *GATA* factors, but *Lyl1* cannot rescue
698 the early *Scl*^{-/-} phenotype. *Blood*. 2007;109(5):1908-16.
- 699 22. Wilkinson AC, Kawata VK, Schutte J, Gao X, Antoniou S, Baumann C, et al. Single-cell analyses
700 of regulatory network perturbations using enhancer-targeting TALEs suggest novel roles for *PU.1*
701 during haematopoietic specification. *Development*. 2014;141(20):4018-30.
- 702 23. Nottingham WT, Jarratt A, Burgess M, Speck CL, Cheng JF, Prabhakar S, et al. *Runx1*-mediated
703 hematopoietic stem-cell emergence is controlled by a *Gata/Ets/SCL*-regulated enhancer. *Blood*.
704 2007;110(13):4188-97.
- 705 24. Gottgens B, Broccardo C, Sanchez MJ, Deveaux S, Murphy G, Gothert JR, et al. The *scl* +18/19
706 stem cell enhancer is not required for hematopoiesis: identification of a 5' bifunctional hematopoietic-
707 endothelial enhancer bound by *Fli-1* and *Elf-1*. *Molecular and cellular biology*. 2004;24(5):1870-83.
- 708 25. Gottgens B, Nastos A, Kinston S, Piltz S, Delabesse EC, Stanley M, et al. Establishing the
709 transcriptional programme for blood: the *SCL* stem cell enhancer is regulated by a multiprotein
710 complex containing *Ets* and *GATA* factors. *The EMBO journal*. 2002;21(12):3039-50.
- 711 26. Gottgens B, Ferreira R, Sanchez MJ, Ishibashi S, Li J, Spensberger D, et al. *cis*-Regulatory
712 remodeling of the *SCL* locus during vertebrate evolution. *Molecular and cellular biology*.
713 2010;30(24):5741-51.
- 714 27. Calero-Nieto FJ, Ng FS, Wilson NK, Hannah R, Moignard V, Leal-Cervantes AI, et al. Key
715 regulators control distinct transcriptional programmes in blood progenitor and mast cells. *The EMBO*
716 *journal*. 2014;33(11):1212-26.
- 717 28. Hardison RC, Taylor J. Genomic approaches towards finding *cis*-regulatory modules in animals.
718 *Nature reviews Genetics*. 2012;13(7):469-83.
- 719 29. Bee T, Ashley EL, Bickley SR, Jarratt A, Li PS, Sloane-Stanley J, et al. The mouse *Runx1* +23
720 hematopoietic stem cell enhancer confers hematopoietic specificity to both *Runx1* promoters. *Blood*.
721 2009;113(21):5121-4.
- 722 30. Sinclair AM, Gottgens B, Barton LM, Stanley ML, Pardanaud L, Klaine M, et al. Distinct 5' *SCL*
723 enhancers direct transcription to developing brain, spinal cord, and endothelium: neural expression is
724 mediated by *GATA* factor binding sites. *Dev Biol*. 1999;209(1):128-42.
- 725 31. Sanchez M, Gottgens B, Sinclair AM, Stanley M, Begley CG, Hunter S, et al. An *SCL* 3' enhancer
726 targets developing endothelium together with embryonic and adult haematopoietic progenitors.
727 *Development*. 1999;126(17):3891-904.

728 32. Lelieveld SH, Schutte J, Dijkstra MJ, Bawono P, Kinston SJ, Gottgens B, et al. ConBind: motif-
729 aware cross-species alignment for the identification of functional transcription factor binding sites.
730 *Nucleic acids research*. 2015.

731 33. Mikkola HKA, Klintman J, Yang H, Hock H, Schlaeger TM, Fujiwara Y, et al. Haematopoietic
732 stem cells retain long-term repopulating activity and multipotency in the absence of stem-cell
733 leukaemia SCL/tal-1 gene. *Nature*. 2003;421(6922):547-51.

734 34. Hall MA, Curtis DJ, Metcalf D, Elefanty AG, Sourris K, Robb L, et al. The critical regulator of
735 embryonic hematopoiesis, SCL, is vital in the adult for megakaryopoiesis, erythropoiesis, and lineage
736 choice in CFU-S12. *Proceedings of the National Academy of Sciences of the United States of America*.
737 2003;100(3):992-7.

738 35. Capron C, Lecluse Y, Kaushik AL, Foudi A, Lacout C, Sekkai D, et al. The SCL relative LYL-1 is
739 required for fetal and adult hematopoietic stem cell function and B-cell differentiation. *Blood*.
740 2006;107(12):4678-86.

741 36. Souroullas GP, Salmon JM, Sablitzky F, Curtis DJ, Goodell MA. Adult hematopoietic stem and
742 progenitor cells require either Lyl1 or Scl for survival. *Cell stem cell*. 2009;4(2):180-6.

743 37. Licht JD. AML1 and the AML1-ETO fusion protein in the pathogenesis of t(8;21) AML.
744 *Oncogene*. 2001;20(40):5660-79.

745 38. Tijssen MR, Cvejic A, Joshi A, Hannah RL, Ferreira R, Forrai A, et al. Genome-wide analysis of
746 simultaneous GATA1/2, RUNX1, FLI1, and SCL binding in megakaryocytes identifies hematopoietic
747 regulators. *Developmental cell*. 2011;20(5):597-609.

748 39. Shivdasani RA, Mayer EL, Orkin SH. Absence of blood formation in mice lacking the T-cell
749 leukaemia oncoprotein tal-1/SCL. *Nature*. 1995;373(6513):432-4.

750 40. Taniuchi I, Sunshine MJ, Festenstein R, Littman DR. Evidence for distinct CD4 silencer functions
751 at different stages of thymocyte differentiation. *Molecular cell*. 2002;10(5):1083-96.

752 41. Taniuchi I, Osato M, Egawa T, Sunshine MJ, Bae SC, Komori T, et al. Differential requirements
753 for Runx proteins in CD4 repression and epigenetic silencing during T lymphocyte development. *Cell*.
754 2002;111(5):621-33.

755 42. Takahashi K, Yamanaka S. Induction of pluripotent stem cells from mouse embryonic and adult
756 fibroblast cultures by defined factors. *Cell*. 2006;126(4):663-76.

757 43. Graf T, Enver T. Forcing cells to change lineages. *Nature*. 2009;462(7273):587-94.

758 44. Batta K, Florkowska M, Kouskoff V, Lacaud G. Direct reprogramming of murine fibroblasts to
759 hematopoietic progenitor cells. *Cell reports*. 2014;9(5):1871-84.

760 45. Riddell J, Gazit R, Garrison BS, Guo G, Saadatpour A, Mandal PK, et al. Reprogramming
761 committed murine blood cells to induced hematopoietic stem cells with defined factors. *Cell*.
762 2014;157(3):549-64.

763 46. Peter IS, Davidson EH. A gene regulatory network controlling the embryonic specification of
764 endoderm. *Nature*. 2011;474(7353):635-9.

765 47. Pinto do OP, Kolterud A, Carlsson L. Expression of the LIM-homeobox gene LH2 generates
766 immortalized steel factor-dependent multipotent hematopoietic precursors. *The EMBO journal*.
767 1998;17(19):5744-56.

768 48. Dexter TM, Allen TD, Scott D, Teich NM. Isolation and characterisation of a bipotential
769 haematopoietic cell line. *Nature*. 1979;277(5696):471-4.

770 49. Busch K, Klapproth K, Barile M, Flossdorf M, Holland-Letz T, Schlenner SM, et al. Fundamental
771 properties of unperturbed haematopoiesis from stem cells in vivo. *Nature*. 2015;518(7540):542-6.

772 50. Langmead B, Salzberg SL. Fast gapped-read alignment with Bowtie 2. *Nat Methods*.
773 2012;9(4):357-9.

774 51. Zhang Y, Liu T, Meyer CA, Eeckhoute J, Johnson DS, Bernstein BE, et al. Model-based analysis
775 of ChIP-Seq (MACS). *Genome Biol*. 2008;9(9):R137.

776 52. Sanchez-Castillo M, Ruau D, Wilkinson AC, Ng FS, Hannah R, Diamanti E, et al. CODEX: a next-
777 generation sequencing experiment database for the haematopoietic and embryonic stem cell
778 communities. *Nucleic acids research*. 2015;43(Database issue):D1117-23.

- 779 53. Xu W, Kee BL. Growth factor independent 1B (Gfi1b) is an E2A target gene that modulates
780 Gata3 in T-cell lymphomas. *Blood*. 2007;109(10):4406-14.
- 781 54. Wang W, Yang J, Liu H, Lu D, Chen X, Zenonos Z, et al. Rapid and efficient reprogramming of
782 somatic cells to induced pluripotent stem cells by retinoic acid receptor gamma and liver receptor
783 homolog 1. *Proceedings of the National Academy of Sciences of the United States of America*.
784 2011;108(45):18283-8.
- 785 55. Koller D, Friedman N. *Probabilistic Graphical Models: Principles and Techniques*. 2009.
- 786 56. Longabaugh WJ, Davidson EH, Bolouri H. Computational representation of developmental
787 genetic regulatory networks. *Dev Biol*. 2005;283(1):1-16.

788

789

790 **Figure legends**

791 **Figure 1: Identification of haematopoietic active *cis*-regulatory regions.** (a) UCSC
792 screenshot of the *Erg* gene locus for ChIP-Sequencing data for nine haematopoietic TFs (ERG,
793 FLI1, GATA2, GFI1B, LYL1, MEIS1, PU.1, RUNX1 and TAL1 (15)) and for H3K27ac (27)
794 in HPC7 cells. Highlighted are all regions of the *Erg* gene locus that are acetylated at H3K27
795 and are bound by three or more TFs. Numbers indicate the distance (in kb) from the ATG start
796 codon. (b) Summary of the identification of candidate *cis*-regulatory regions for all nine TFs
797 and subsequent analysis in transgenic mouse assays. The inspection of the nine gene loci and
798 the application of the selection criteria (≥ 3 TFs bound and H3K27ac) identified a total of 49
799 candidate *cis*-regulatory regions. The heatmap shows the binding pattern of the nine TFs to all
800 candidate regulatory elements in HPC7 cells: green = bound, grey = unbound. Haematopoietic
801 activity in E11.5 transgenic mice is indicated by the font color: black = active, red = not active.
802 Grey indicates genomic repeat regions that were not tested in transgenic mice. Detailed
803 experimental data corresponding to the summary heatmap can be found in Fig. 1-figure
804 supplement 1-8. (c) Haematopoietic activity of the five candidate *Erg cis*-regulatory regions
805 was determined in E11.5 transgenic mouse assays. Shown are X-Gal-stained whole-mount
806 embryos and paraffin sections of the dorsal aorta (DA, ventral side on the left/top) and foetal
807 liver (FL), sites of definitive haematopoiesis. Colour coding as in B.

808

809 **Figure 1 – figure supplement 1: Identification of haematopoietic active *cis*-regulatory**
810 **elements for *Fli1*.** (a) The candidate *cis*-regulatory elements were identified by ChIP-Seq
811 analysis of the TFs ERG, FLI1, GATA2, GFI1B, LYL1, MEIS1, PU.1, RUNX1 and TAL1 as
812 well as H3K27 acetylation in the haematopoietic stem/progenitor cell line HPC7. Highlighted
813 in pink are the candidate *cis*-regulatory regions which are bound by at least three of the nine
814 TFs and showed H3K27 acetylation. The numbering represents the direction and distance in
815 kilobases from the start codon ATG (pro = promoter). (b) Candidate regions were assayed for
816 haematopoietic enhancer activity in mouse transient transgenic embryos. X-Gal stained whole-
817 mount E11.5 embryos and paraffin sections of the dorsal aorta (DA; longitudinal section,
818 ventral side on the left/top) and foetal liver (FL) are shown for the candidate *cis*-regulatory
819 regions. Transgenic mouse data are not shown for previously published regions, but relevant
820 publications are listed.

821

822 **Figure 1 – figure supplement 2: Identification of haematopoietic active *cis*-regulatory**
823 **elements for *Gata2*.** (a) The candidate *cis*-regulatory elements were identified by ChIP-Seq
824 analysis of the TFs ERG, FLI1, GATA2, GFI1B, LYL1, MEIS1, PU.1, RUNX1 and TAL1 as
825 well as H3K27 acetylation in the haematopoietic stem/progenitor cell line HPC7. Highlighted
826 in pink are the candidate *cis*-regulatory regions which are bound by at least three of the nine
827 TFs and showed H3K27 acetylation. The numbering represents the direction and distance in
828 kilobases from the start codon ATG (pro = promoter). (b) Candidate regions were assayed for
829 haematopoietic enhancer activity in mouse transient transgenic embryos. X-Gal stained whole-
830 mount E11.5 embryos and paraffin sections of the dorsal aorta (DA; longitudinal section,
831 ventral side on the left/top) and foetal liver (FL) are shown for the candidate *cis*-regulatory

832 regions. Transgenic mouse data are not shown for previously published regions, but relevant
833 publications are listed.

834

835 **Figure 1 – figure supplement 3: Identification of haematopoietic active *cis*-regulatory**
836 **elements for *Gfi1b*.** (a) The candidate *cis*-regulatory elements were identified by ChIP-Seq
837 analysis of the TFs ERG, FLI1, GATA2, GFI1B, LYL1, MEIS1, PU.1, RUNX1 and TAL1 as
838 well as H3K27 acetylation in the haematopoietic stem/progenitor cell line HPC7. Highlighted
839 in pink are the candidate *cis*-regulatory regions which are bound by at least three of the nine
840 TFs and showed H3K27 acetylation. The numbering represents the direction and distance in
841 kilobases from the start codon ATG (pro = promoter). (b) All candidate regions were
842 previously published regions. Relevant publications are listed.

843

844 **Figure 1 – figure supplement 4: Identification of haematopoietic active *cis*-regulatory**
845 **elements for *Lyl1*.** (a) The candidate *cis*-regulatory elements were identified by ChIP-Seq
846 analysis of the TFs ERG, FLI1, GATA2, GFI1B, LYL1, MEIS1, PU.1, RUNX1 and TAL1 as
847 well as H3K27 acetylation in the haematopoietic stem/progenitor cell line HPC7. Highlighted
848 in pink are the candidate *cis*-regulatory regions which are bound by at least three of the nine
849 TFs and showed H3K27 acetylation. The numbering represents the direction and distance in
850 kilobases from the start codon ATG (pro = promoter). (b) Candidate regions were assayed for
851 haematopoietic enhancer activity in mouse transient transgenic embryos. X-Gal stained whole-
852 mount E11.5 embryos and paraffin sections of the dorsal aorta (DA; longitudinal section,
853 ventral side on the left/top) and foetal liver (FL) are shown for the candidate *cis*-regulatory
854 regions. Transgenic mouse data are not shown for previously published regions, but relevant
855 publications are listed.

856

857 **Figure 1 – figure supplement 5: Identification of haematopoietic active *cis*-regulatory**
858 **elements for *Meis1*.** (a) The candidate *cis*-regulatory elements were identified by ChIP-Seq
859 analysis of the TFs ERG, FLI1, GATA2, GFI1B, LYL1, MEIS1, PU.1, RUNX1 and TAL1 as
860 well as H3K27 acetylation in the haematopoietic stem/progenitor cell line HPC7. Highlighted
861 in pink are the candidate *cis*-regulatory regions which are bound by at least three of the nine
862 TFs and showed H3K27 acetylation. The numbering represents the direction and distance in
863 kilobases from the start codon ATG (pro = promoter). (b) Candidate regions were assayed for
864 haematopoietic enhancer activity in mouse transient transgenic embryos. X-Gal stained whole-
865 mount E11.5 embryos and paraffin sections of the dorsal aorta (DA; longitudinal section,
866 ventral side on the left/top) and foetal liver (FL) are shown for the candidate *cis*-regulatory
867 regions.

868

869 **Figure 1 – figure supplement 6: Identification of haematopoietic active *cis*-regulatory**
870 **elements for *Runx1*.** (a) The candidate *cis*-regulatory elements were identified by ChIP-Seq
871 analysis of the TFs ERG, FLI1, GATA2, GFI1B, LYL1, MEIS1, PU.1, RUNX1 and TAL1 as
872 well as H3K27 acetylation in the haematopoietic stem/progenitor cell line HPC7. Highlighted
873 in pink are the candidate *cis*-regulatory regions which are bound by at least three of the nine
874 TFs and showed H3K27 acetylation. The numbering represents the direction and distance in
875 kilobases from the start codon ATG (pro = promoter). (b) E10 embryos and cryosections of the
876 DA (transverse; ventral down) and FL are shown. For the *Runx1*+204 region, a larger 12 kb
877 fragment (chr16:92,620,915-92,631,936, mm9) was used for transient transgenesis, but similar
878 results were obtained with the +204 fragment alone (data not shown). The +24 element was
879 tested in conjunction with the +23 and did not change its tissue specificity (Bee et al., 2010).

880 Preliminary data show that the +24 on its own does not mediate robust tissue specific
881 expression of reporter genes. Transgenic mouse data are not shown for previously published
882 regions, but relevant publications are listed.

883

884 **Figure 1 – figure supplement 7: Identification of haematopoietic active *cis*-regulatory**
885 **elements for *Spi1*.** (a) The candidate *cis*-regulatory elements were identified by ChIP-Seq
886 analysis of the TFs ERG, FLI1, GATA2, GFI1B, LYL1, MEIS1, PU.1, RUNX1 and TAL1 as
887 well as H3K27 acetylation in the haematopoietic stem/progenitor cell line HPC7. Highlighted
888 in pink are the candidate *cis*-regulatory regions which are bound by at least three of the nine
889 TFs and showed H3K27 acetylation. The numbering represents the direction and distance in
890 kilobases from the start codon ATG (pro = promoter). (b) All candidate regions were
891 previously published regions. Relevant publications are listed.

892

893 **Figure 1 – figure supplement 8: Identification of haematopoietic active *cis*-regulatory**
894 **elements for *Tal1*.** (a) The candidate *cis*-regulatory elements were identified by ChIP-Seq
895 analysis of the TFs ERG, FLI1, GATA2, GFI1B, LYL1, MEIS1, PU.1, RUNX1 and TAL1 as
896 well as H3K27 acetylation in the haematopoietic stem/progenitor cell line HPC7. Highlighted
897 in pink are the candidate *cis*-regulatory regions which are bound by at least three of the nine
898 TFs and showed H3K27 acetylation. The numbering is based on the distance (in kb) to
899 promoter 1a. (b) All candidate regions were previously published regions. Relevant
900 publications are listed.

901

902 **Figure 1 – figure supplement 9: Number of PCR and LacZ positive transgenic embryos**
903 **(E10.5-11.5) for each regulatory region.**

904

905 **Figure 2: Comparison of TF binding pattern at haematopoietic active *cis*-regulatory**
906 **regions in two haematopoietic progenitor cell lines, HPC7 and 416b. (a)** UCSC screenshot
907 of the *Erg* gene locus for ChIP-Sequencing data for nine haematopoietic TFs (ERG, FLI1,
908 GATA2, GFI1B, LYL1, MEIS1, PU.1, RUNX1 and TAL1) and for H3K27ac in 416b cells.
909 Highlighted are those haematopoietic active *Erg cis*-regulatory regions that were identified
910 based on acetylation of H3K27 and TF binding in HPC7 cells followed by transgenic mouse
911 assays. Numbers indicate the distance (in kb) from the ATG start codon. **(b)** Hierarchical
912 clustering of the binding profiles for HPC7, 416b and other published datasets. The heatmap
913 shows the pairwise correlation coefficient of peak coverage data between pairs of samples in
914 the row and column. The order of the samples is identical in columns and rows. Details about
915 samples listed can be found in Fig. 2-figure supplement 9. **(c)** Pair-wise analysis of binding of
916 the nine TFs to haematopoietic active *cis*-regulatory regions of the nine TFs in HPC7 versus
917 416b cells. Green = bound in both cells types, blue = only bound in 416b cells, orange = only
918 bound in HPC7 cells, grey = not bound in either cell type.

919

920 **Figure 2 – figure supplement 1: UCSC screenshot for the *Fli1* gene locus demonstrating**
921 **binding patterns for nine key haematopoietic TFs and H3K27ac in 416b cells.** Highlighted
922 in pink are *cis*-regulatory regions that were identified based on the selection criteria (≥ 3 TFs
923 bound and H3K27ac) in HPC7 cells and were shown to possess haematopoietic activity. The
924 numbering represents the distance (in kb) from the start codon ATG.

925

926 **Figure 2 – figure supplement 2: UCSC screenshot for the *Gata2* gene locus demonstrating**
927 **binding patterns for nine key haematopoietic TFs and H3K27ac in 416b cells.** Highlighted

928 in pink are *cis*-regulatory regions that were identified based on the selection criteria (≥ 3 TFs
929 bound and H3K27ac) in HPC7 cells and were shown to possess haematopoietic activity. The
930 numbering represents the distance (in kb) from the start codon ATG.

931

932 **Figure 2 – figure supplement 3: UCSC screenshot for the *Gfi1b* gene locus demonstrating**
933 **binding patterns for nine key haematopoietic TFs and H3K27ac in 416b cells.** Highlighted
934 in pink are *cis*-regulatory regions that were identified based on the selection criteria (≥ 3 TFs
935 bound and H3K27ac) in HPC7 cells and were shown to possess haematopoietic activity. The
936 numbering represents the distance (in kb) from the start codon ATG.

937

938 **Figure 2 – figure supplement 4: UCSC screenshot for the *Lyl1* gene locus demonstrating**
939 **binding patterns for nine key haematopoietic TFs and H3K27ac in 416b cells.** Highlighted
940 in pink is the promoter (“pro”) that was identified based on the selection criteria (≥ 3 TFs
941 bound and H3K27ac) in HPC7 cells and was shown to possess haematopoietic activity. The
942 promoter is labelled with “pro”.

943

944 **Figure 2 – figure supplement 5: UCSC screenshot for the *Meis1* gene locus demonstrating**
945 **binding patterns for nine key haematopoietic TFs and H3K27ac in 416b cells.** Highlighted
946 in pink is the *cis*-regulatory region that was identified based on the selection criteria (≥ 3 TFs
947 bound and H3K27ac) in HPC7 cells and was shown to possess haematopoietic activity. The
948 numbering represents the distance (in kb) from the start codon ATG.

949

950 **Figure 2 – figure supplement 6: UCSC screenshot for the *Runx1* gene locus**
951 **demonstrating binding patterns for nine key haematopoietic TFs and H3K27ac in 416b**
952 **cells.** Highlighted in pink are *cis*-regulatory regions that were identified based on the selection
953 criteria (≥ 3 TFs bound and H3K27ac) in HPC7 cells and were subsequently shown to possess
954 haematopoietic activity. The numbering represents the distance (in kb) from the start codon
955 ATG.

956

957 **Figure 2 – figure supplement 7: UCSC screenshot for the *Spi1* gene locus demonstrating**
958 **binding patterns for nine key haematopoietic TFs and H3K27ac in 416b cells.** Highlighted
959 in pink is the *cis*-regulatory region that was identified based on the selection criteria (≥ 3 TFs
960 bound and H3K27ac) in HPC7 cells and was shown to possess haematopoietic activity. The
961 numbering represents the distance (in kb) from the start codon ATG.

962

963 **Figure 2 – figure supplement 8: UCSC screenshot for the *Tall* gene locus demonstrating**
964 **binding patterns for nine key haematopoietic TFs and H3K27ac in 416b cells.** Highlighted
965 in pink are *cis*-regulatory regions that were identified based on the selection criteria (≥ 3 TFs
966 bound and H3K27ac) in HPC7 cells and were shown to possess haematopoietic activity. The
967 numbering the distance (in kb) from promoter 1a.

968

969 **Figure 2 – figure supplement 9: List of ChIP-Seq samples included in the heatmap in**
970 **Figure 2b.**

971

972 **Figure 3: TFBS mutagenesis reveals enhancer-dependent effects of TF binding on gene**
973 **expression. (a)** Multiple species alignment of mouse (mm9), human (hg19), dog (canFam2),
974 opossum (monDom5) and platypus (ornAna1) sequences for the *Erg*+65 region. Nucleotides
975 highlighted in black are conserved between all species analysed, nucleotides highlighted in
976 grey are conserved between four of five species. Transcription factor binding sites (TFBS) are
977 highlighted in: blue = Ebox, purple = Ets, green = Gata, yellow = Gfi, red = Meis. The
978 nucleotides changed to mutate the TFBSs are indicated below the alignment. All binding sites
979 of one motif family (e.g. all Ebox motifs) were mutated simultaneously. **(b)** Luciferase assay
980 for the *Erg*+65 wild-type and mutant enhancer in stably transfected 416b cells. Each bar
981 represents the averages of at least three independent experiments with three to four replicates
982 within each experiment. Results are shown relative to the wild-type enhancer activity, which is
983 set to 100%. Error bars represent the standard error of the mean (SEM). Stars indicate
984 significance: ** = p-value < 0.01, *** = p-value < 0.001. P-values were calculated using t-
985 tests, followed by the Fisher's method. **(c)** Summary of luciferase assay results for all 19 high-
986 confidence haematopoietic active regulatory regions. Relative luciferase activity is illustrated
987 in shades of blue (down-regulation) and red (up-regulation). Crossed-out grey boxes indicate
988 that there is no motif for the TF and/or the TF does not bind to the region. Detailed results and
989 corresponding alignments with highlighted TFBSs and their mutations can be found in Figure
990 3-figure supplements 1-18.

991

992 **Figure 3 – figure supplement 1: Multiple species alignment and luciferase assay results**
993 **for *Erg*+75. (a)** Multiple species alignment (MSA) with the following species: mouse (mm9),
994 human (hg19), dog (canFam2), opossum (monDom5) and platypus (ornAna1). Nucleotides
995 highlighted in black are conserved between all species analysed, nucleotides highlighted in
996 grey are conserved between four of five species. Transcription factor binding sites (TFBS) are

997 highlighted in: blue = Ebox, purple = Ets, yellow = Gfi. The nucleotides that were changed to
998 mutate the TFBSs are indicated below the MSA. All conserved binding sites of one motif
999 family (e.g. all Ebox motifs) were mutated simultaneously. Where TF binding was observed in
1000 ChIP-Seq experiments in 416b cells, but the TFBS was not conserved, the motifs present in the
1001 mouse sequence only were mutated. **(b)** For the luciferase reporter assays in stably transfected
1002 416b cells the averages of at least three independent experiments with three to four replicates
1003 within each experiment are shown. Error bars represent the standard error of the mean (SEM).
1004 Stars indicate significance: ** = p-value < 0.01, *** = p-value < 0.001. P-values were
1005 generated using t-tests, followed by the Fisher's method and if necessary Stouffer's z trend.

1006

1007 **Figure 3 – figure supplement 2: Multiple species alignment and luciferase assay results**
1008 **for *Erg+85*.** **(a)** Multiple species alignment (MSA) with the following species: mouse (mm9),
1009 human (hg19), dog (canFam2), opossum (monDom5) and platypus (ornAna1). Nucleotides
1010 highlighted in black are conserved between all species analysed, nucleotides highlighted in
1011 grey are conserved between four of five species. Transcription factor binding sites (TFBS) are
1012 highlighted in: blue = Ebox, purple = Ets, green = Gata, yellow = Gfi. The nucleotides that
1013 were changed to mutate the TFBSs are indicated below the MSA. All conserved binding sites
1014 of one motif family (e.g. all Ebox motifs) were mutated simultaneously. **(b)** For the luciferase
1015 reporter assays in stably transfected 416b cells the averages of at least three independent
1016 experiments with three to four replicates within each experiment are shown. Error bars
1017 represent the standard error of the mean (SEM). Stars indicate significance: *** =
1018 p-value < 0.001. P-values were generated using t-tests, followed by the Fisher's method and if
1019 necessary Stouffer's z trend.

1020

1021 **Figure 3 – figure supplement 3: Multiple species alignment and luciferase assay results**
1022 **for *Fli1+12*.** (a) Multiple species alignment (MSA) with the following species: mouse (mm9),
1023 human (hg19), dog (canFam2), opossum (monDom5) and platypus (ornAna1). Nucleotides
1024 highlighted in black are conserved between all species analysed, nucleotides highlighted in
1025 grey are conserved between four of five species. Transcription factor binding sites (TFBS) are
1026 highlighted in: blue = Ebox, purple = Ets. The nucleotides that were changed to mutate the
1027 TFBSs are indicated below the MSA. All conserved binding sites of one motif family (e.g. all
1028 Ebox motifs) were mutated simultaneously. (b) For the luciferase reporter assays in stably
1029 transfected 416b cells the averages of at least three independent experiments with three to four
1030 replicates within each experiment are shown. Error bars represent the standard error of the
1031 mean (SEM). Stars indicate significance: *** = p-value < 0.001. P-values were generated using
1032 t-tests, followed by the Fisher's method and if necessary Stouffer's z trend.

1033

1034 **Figure 3 – figure supplement 4: Multiple species alignment and luciferase assay results**
1035 **for *Gata2-93*.** (a) Multiple species alignment (MSA) with the following species: mouse
1036 (mm9), human (hg19), dog (canFam2), opossum (monDom5) and platypus (ornAna1).
1037 Nucleotides highlighted in black are conserved between all species analysed, nucleotides
1038 highlighted in grey are conserved between four of five species. Transcription factor binding
1039 sites (TFBS) are highlighted in: blue = Ebox, purple = Ets, green = Gata, red = Meis, turquoise
1040 = Runt. The nucleotides that were changed to mutate the TFBSs are indicated below the MSA.
1041 All conserved binding sites of one motif family (e.g. all Ebox motifs) were mutated
1042 simultaneously. (b) For the luciferase reporter assays in stably transfected 416b cells the
1043 averages of at least three independent experiments with three to four replicates within each
1044 experiment are shown. Error bars represent the standard error of the mean (SEM). Stars

1045 indicate significance: ** = p-value < 0.01, *** = p-value < 0.001. P-values were generated
1046 using t-tests, followed by the Fisher's method and if necessary Stouffer's z trend.

1047

1048 **Figure 3 – figure supplement 5: Multiple species alignment and luciferase assay results**

1049 **for *Gata2+3*. (a)** Multiple species alignment (MSA) with the following species: mouse (mm9),

1050 human (hg19), dog (canFam2), opossum (monDom5) and platypus (ornAna1). Nucleotides

1051 highlighted in black are conserved between all species analysed, nucleotides highlighted in

1052 grey are conserved between four of five species. Transcription factor binding sites (TFBS) are

1053 highlighted in: blue = Ebox, purple = Ets, green = Gata. The nucleotides that were changed to

1054 mutate the TFBSs are indicated below the MSA. All conserved binding sites of one motif

1055 family (e.g. all Ebox motifs) were mutated simultaneously. **(b)** For the luciferase reporter

1056 assays in stably transfected 416b cells the averages of at least three independent experiments

1057 with three to four replicates within each experiment are shown. Error bars represent the

1058 standard error of the mean (SEM). Stars indicate significance: *** = p-value < 0.001. P-values

1059 were generated using t-tests, followed by the Fisher's method and if necessary Stouffer's z

1060 trend.

1061

1062 **Figure 3 – figure supplement 6: Multiple species alignment and luciferase assay results**

1063 **for *Gfi1b+16*. (a)** Multiple species alignment (MSA) with the following species: mouse

1064 (mm9), human (hg19), dog (canFam2), opossum (monDom5) and platypus (ornAna1).

1065 Nucleotides highlighted in black are conserved between all species analysed, nucleotides

1066 highlighted in grey are conserved between four of five species. Transcription factor binding

1067 sites (TFBS) are highlighted in: blue = Ebox, purple = Ets, green = Gata, yellow = Gfi, red =

1068 Meis, turquoise = Runt. The nucleotides that were changed to mutate the TFBSs are indicated

1069 below the MSA. All conserved binding sites of one motif family (e.g. all Ebox motifs) were
1070 mutated simultaneously. Where TF binding was observed in ChIP-Seq experiments in 416b
1071 cells, but the TFBS was not conserved, the motifs present in the mouse sequence only were
1072 mutated. **(b)** For the luciferase reporter assays in stably transfected 416b cells the averages of
1073 at least three independent experiments with three to four replicates within each experiment are
1074 shown. Error bars represent the standard error of the mean (SEM). Stars indicate significance:
1075 ** = p-value < 0.01, *** = p-value < 0.001. P-values were generated using t-tests, followed by
1076 the Fisher's method and if necessary Stouffer's z trend.

1077

1078 **Figure 3 – figure supplement 7: Multiple species alignment and luciferase assay results**
1079 **for *Gfi1b*+17. (a)** Multiple species alignment (MSA) with the following species: mouse
1080 (mm9), human (hg19), dog (canFam2), opossum (monDom5) and platypus (ornAna1).
1081 Nucleotides highlighted in black are conserved between all species analysed, nucleotides
1082 highlighted in grey are conserved between four of five species. Transcription factor binding
1083 sites (TFBS) are highlighted in: blue = Ebox, purple = Ets, green = Gata, yellow = Gfi, red =
1084 Meis. The nucleotides that were changed to mutate the TFBSs are indicated below the MSA.
1085 All conserved binding sites of one motif family (e.g. all Ebox motifs) were mutated
1086 simultaneously. **(b)** For the luciferase reporter assays in stably transfected 416b cells the
1087 averages of at least three independent experiments with three to four replicates within each
1088 experiment are shown. Error bars represent the standard error of the mean (SEM). Stars
1089 indicate significance: ** = p-value < 0.01, *** = p-value < 0.001. P-values were generated
1090 using t-tests, followed by the Fisher's method and if necessary Stouffer's z trend.

1091

1092 **Figure 3 – figure supplement 8: Multiple species alignment and luciferase assay results**
1093 **for *Lyf1* promoter. (a)** Multiple species alignment (MSA) with the following species: mouse
1094 (mm9), human (hg19), dog (canFam2) and opossum (monDom5). Nucleotides highlighted in
1095 black are conserved between all species analysed, nucleotides highlighted in grey are
1096 conserved between three of four species. Transcription factor binding sites (TFBS) are
1097 highlighted in: purple = Ets, green = Gata. The nucleotides that were changed to mutate the
1098 TFBSs are indicated below the MSA. All conserved binding sites of one motif family (e.g. all
1099 Ets motifs) were mutated simultaneously. **(b)** For the luciferase reporter assays in stably
1100 transfected 416b cells the averages of at least three independent experiments with three to four
1101 replicates within each experiment are shown. Error bars represent the standard error of the
1102 mean (SEM). Stars indicate significance: *** = p-value < 0.001. P-values were generated using
1103 t-tests, followed by the Fisher's method and if necessary Stouffer's z trend.

1104

1105 **Figure 3 – figure supplement 9: Multiple species alignment and luciferase assay results**
1106 **for *Meis1+48*. (a)** Multiple species alignment (MSA) with the following species: mouse
1107 (mm9), human (hg19), dog (canFam2), opossum (monDom5) and platypus (ornAna1).
1108 Nucleotides highlighted in black are conserved between all species analysed, nucleotides
1109 highlighted in grey are conserved between four of five species. Transcription factor binding
1110 sites (TFBS) are highlighted in: purple = Ets, green = Gata, yellow = Gfi, red = Meis. The
1111 nucleotides that were changed to mutate the TFBSs are indicated below the MSA. All
1112 conserved binding sites of one motif family (e.g. all Ets motifs) were mutated simultaneously.
1113 **(b)** For the luciferase reporter assays in stably transfected 416b cells the averages of at least
1114 three independent experiments with three to four replicates within each experiment are shown.
1115 Error bars represent the standard error of the mean (SEM). Stars indicate significance: *** =

1116 p-value < 0.001. P-values were generated using t-tests, followed by the Fisher's method and if
1117 necessary Stouffer's z trend.

1118

1119 **Figure 3 – figure supplement 10: Multiple species alignment and luciferase assay results**
1120 **for *Spi1-14*.** (a) Multiple species alignment (MSA) with the following species: mouse (mm9),
1121 human (hg19), dog (canFam2), opossum (monDom5) and platypus (ornAna1). Nucleotides
1122 highlighted in black are conserved between all species analysed, nucleotides highlighted in
1123 grey are conserved between four of five species. Transcription factor binding sites (TFBS) are
1124 highlighted in: blue = Ebox, purple = Ets, turquoise = Runt. The nucleotides that were changed
1125 to mutate the TFBSs are indicated below the MSA. All conserved binding sites of one motif
1126 family (e.g. all Ebox motifs) were mutated simultaneously. (b) For the luciferase reporter
1127 assays in stably transfected 416b cells the averages of at least three independent experiments
1128 with three to four replicates within each experiment are shown. Error bars represent the
1129 standard error of the mean (SEM). Stars indicate significance: ** = p-value < 0.01, *** =
1130 p-value < 0.001. P-values were generated using t-tests, followed by the Fisher's method and if
1131 necessary Stouffer's z trend.

1132

1133 **Figure 3 – figure supplement 11: Multiple species alignment and luciferase assay results**
1134 **for *Runx1-59*.** (a) Multiple species alignment (MSA) with the following species: mouse
1135 (mm9), human (hg19) and dog (canFam2). Nucleotides highlighted in black are conserved
1136 between all species analysed, nucleotides highlighted in grey are conserved between two of
1137 three species. Transcription factor binding sites (TFBS) are highlighted in: blue = Ebox, purple
1138 = Ets, green = Gata, red = Meis. The nucleotides that were changed to mutate the TFBSs are
1139 indicated below the MSA. All conserved binding sites of one motif family (e.g. all Ebox

1140 motifs) were mutated simultaneously. **(b)** For the luciferase reporter assays in stably
1141 transfected 416b cells the averages of at least three independent experiments with three to four
1142 replicates within each experiment are shown. Error bars represent the standard error of the
1143 mean (SEM). Stars indicate significance: *** = p-value < 0.001. P-values were generated using
1144 t-tests, followed by the Fisher's method and if necessary Stouffer's z trend.

1145

1146 **Figure 3 – figure supplement 12: Multiple species alignment and luciferase assay results**

1147 **for *Runx1+3*.** **(a)** Multiple species alignment (MSA) with the following species: mouse
1148 (mm9), human (hg19), dog (canFam2), opossum (monDom5) and platypus (ornAna1).
1149 Nucleotides highlighted in black are conserved between all species analysed, nucleotides
1150 highlighted in grey are conserved between four of five species. Transcription factor binding
1151 sites (TFBS) are highlighted in: blue = Ebox, purple = Ets, green = Gata, yellow = Gfi, red =
1152 Meis, turquoise = Runt. The nucleotides that were changed to mutate the TFBSs are indicated
1153 below the MSA. All conserved binding sites of one motif family (e.g. all Ets motifs) were
1154 mutated simultaneously. Where TF binding was observed in ChIP-Seq experiments in 416b
1155 cells, but the TFBS was not conserved, the motifs present in the mouse sequence only were
1156 mutated. **(b)** For the luciferase reporter assays in stably transfected 416b cells the averages of
1157 at least three independent experiments with three to four replicates within each experiment are
1158 shown. Error bars represent the standard error of the mean (SEM). Stars indicate significance:
1159 * = p-value < 0.05, ** = p-value < 0.01, *** = p-value < 0.001. P-values were generated using
1160 t-tests, followed by the Fisher's method and if necessary Stouffer's z trend.

1161

1162 **Figure 3 – figure supplement 13: Multiple species alignment and luciferase assay results**

1163 **for *Runx1+23*.** **(a)** Multiple species alignment (MSA) with the following species: mouse

1164 (mm9), human (hg19), dog (canFam2) and opossum (monDom5). Nucleotides highlighted in
1165 black are conserved between all species analysed, nucleotides highlighted in grey are
1166 conserved between three to four species. Transcription factor binding sites (TFBS) are
1167 highlighted in: blue = Ebox, purple = Ets, green = Gata, red = Meis, turquoise = Runt. The
1168 nucleotides that were changed to mutate the TFBSs are indicated below the MSA. All
1169 conserved binding sites of one motif family (e.g. all Ebox motifs) were mutated
1170 simultaneously. **(b)** For the luciferase reporter assays in stably transfected 416b cells the
1171 averages of at least three independent experiments with three to four replicates within each
1172 experiment are shown. Error bars represent the standard error of the mean (SEM). Stars
1173 indicate significance: * = p-value < 0.05, *** = p-value < 0.001. P-values were generated using
1174 t-tests, followed by the Fisher's method and if necessary Stouffer's z trend.

1175

1176 **Figure 3 – figure supplement 14: Multiple species alignment and luciferase assay results**
1177 **for *Runx1*+110.** **(a)** Multiple species alignment (MSA) with the following species: mouse
1178 (mm9), human (hg19), dog (canFam2), opossum (monDom5) and platypus (ornAna1).
1179 Nucleotides highlighted in black are conserved between all species analysed, nucleotides
1180 highlighted in grey are conserved between four of five species. Transcription factor binding
1181 sites (TFBS) are highlighted in: blue = Ebox, purple = Ets, green = Gata. The nucleotides that
1182 were changed to mutate the TFBSs are indicated below the MSA. All conserved binding sites
1183 of one motif family (e.g. all Ets motifs) were mutated simultaneously. Where TF binding was
1184 observed in CHIP-Seq experiments in 416b cells, but the TFBS was not conserved, the motifs
1185 present in the mouse sequence only were mutated. **(b)** For the luciferase reporter assays in
1186 stably transfected 416b cells the averages of at least three independent experiments with three
1187 to four replicates within each experiment are shown. Error bars represent the standard error of
1188 the mean (SEM). Stars indicate significance: ** = p-value < 0.01, *** = p-value < 0.001. P-

1189 values were generated using t-tests, followed by the Fisher's method and if necessary
1190 Stouffer's z trend.

1191

1192 **Figure 3 – figure supplement 15: Multiple species alignment and luciferase assay results**

1193 **for *Runx1+204*. (a)** Multiple species alignment (MSA) with the following species: mouse
1194 (mm9), human (hg19), dog (canFam2), opossum (monDom5) and platypus (ornAna1).

1195 Nucleotides highlighted in black are conserved between all species analysed, nucleotides

1196 highlighted in grey are conserved between four of five species. Transcription factor binding

1197 sites (TFBS) are highlighted in: blue = Ebox, purple = Ets, yellow = Gfi, turquoise = Runt. The

1198 nucleotides that were changed to mutate the TFBSs are indicated below the MSA. All

1199 conserved binding sites of one motif family (e.g. all Ets motifs) were mutated simultaneously.

1200 Where TF binding was observed in ChIP-Seq experiments in 416b cells, but the TFBS was not

1201 conserved, the motifs present in the mouse sequence only were mutated. **(b)** For the luciferase

1202 reporter assays in stably transfected 416b cells the averages of at least three independent

1203 experiments with three to four replicates within each experiment are shown. Error bars

1204 represent the standard error of the mean (SEM). Stars indicate significance: *** =

1205 p-value < 0.001. P-values were generated using t-tests, followed by the Fisher's method and if

1206 necessary Stouffer's z trend.

1207

1208 **Figure 3 – figure supplement 16: Multiple species alignment and luciferase assay results**

1209 **for *Tal1-4*. (a)** Multiple species alignment (MSA) with the following species: mouse (mm9),

1210 human (hg19) and dog (canFam2). Nucleotides highlighted in black are conserved between all

1211 species analysed, nucleotides highlighted in grey are conserved between two of three species.

1212 Transcription factor binding sites (TFBS) are highlighted in: purple = Ets. The nucleotides that

1213 were changed to mutate the TFBSs are indicated below the MSA. All conserved binding sites
1214 of the Ets motif family were mutated simultaneously. **(b)** For the luciferase reporter assays in
1215 stably transfected 416b cells the averages of at least three independent experiments with three
1216 to four replicates within each experiment are shown. Error bars represent the standard error of
1217 the mean (SEM). Stars indicate significance: ** = p-value < 0.01. P-values were generated
1218 using t-tests, followed by the Fisher's method and if necessary Stouffer's z trend.

1219

1220 **Figure 3 – figure supplement 17: Multiple species alignment and luciferase assay results**
1221 **for *Tal1+19*.** **(a)** Multiple species alignment (MSA) with the following species: mouse (mm9),
1222 human (hg19), dog (canFam2) and opossum (monDom5). Nucleotides highlighted in black are
1223 conserved between all species analysed, nucleotides highlighted in grey are conserved between
1224 three of four species. Transcription factor binding sites (TFBS) are highlighted in: purple = Ets.
1225 The nucleotides that were changed to mutate the TFBSs are indicated below the MSA. All
1226 conserved binding sites of the Ets motif family were mutated simultaneously. **(b)** For the
1227 luciferase reporter assays in stably transfected 416b cells the averages of at least three
1228 independent experiments with three to four replicates within each experiment are shown. Error
1229 bars represent the standard error of the mean (SEM). Stars indicate significance: *** =
1230 p-value < 0.001. P-values were generated using t-tests, followed by the Fisher's method and if
1231 necessary Stouffer's z trend.

1232

1233 **Figure 3 – figure supplement 18: Multiple species alignment and luciferase assay results**
1234 **for *Tal1+40*.** **(a)** Multiple species alignment (MSA) with the following species: mouse (mm9),
1235 human (hg19) and dog (canFam2). Nucleotides highlighted in black are conserved between all
1236 species analysed, nucleotides highlighted in grey are conserved between two of three species.

1237 Transcription factor binding sites (TFBS) are highlighted in: blue = Ebox, purple = Ets, green =
1238 Gata. The nucleotides that were changed to mutate the TFBSs are indicated below the MSA.
1239 All conserved binding sites of one motif family (e.g. all Ebox motifs) were mutated
1240 simultaneously. **(b)** For the luciferase reporter assays in stably transfected 416b cells the
1241 averages of at least three independent experiments with three to four replicates within each
1242 experiment are shown. Error bars represent the standard error of the mean (SEM). Stars
1243 indicate significance: * = p-value < 0.05. P-values were generated using t-tests, followed by the
1244 Fisher's method and if necessary Stouffer's z trend.

1245

1246 **Figure 3 – figure supplement 19: List of TF binding sites and the TFs that bind to them.**

1247

1248 **Figure 3 – figure supplement 20: List of co-ordinates and primer sequences for the**
1249 **regulatory regions analysed in this study.**

1250

1251 **Figure 4: A three-tier DBN incorporating transcriptional regulatory information can**
1252 **recapitulate the HSPC expression state. (a)** Representation of the complete network diagram
1253 generated using the Biotapestry software (56). **(b)** Schematic diagram describing the DBN which
1254 contains three tiers: I. TF binding motifs within regulatory regions, II. *cis*-regulatory regions
1255 influencing the expression levels of the various TFs, and III. genes encoding the TFs. The
1256 output of tier III, namely the expression levels of the TF, feed back into the TF binding at the
1257 various motifs of tier I. The model therefore is comprised of successive time slices (t). **(c)**
1258 Simulation of a single cell over time. The expression levels of all 9 TFs are the same at the
1259 beginning (0.5). The simulation rapidly stabilizes with characteristic TF expression levels. **(d)**

1260 Simulation of a cell population by running the model 1000 times. The scale of the x-axis is
1261 linear. Each simulation was run as described in (c).

1262

1263 **Figure 4 - figure supplement 1: Simulation of a single cell over time with different**
1264 **expression levels at the beginning.** The simulation rapidly stabilizes with characteristic TF
1265 expression levels irrespective of the starting conditions. **(a)** The expression levels of all 9 TFs
1266 are 0.2 at the start of the simulation. **(b)** The expression levels of all 9 TFs are 0.8 at the start of
1267 the simulation. **(c)** The expression levels for FLI1, RUNX1 and TAL1 are set to be 0.5 at the
1268 beginning, with all other TFs not being expressed (value of 0).

1269

1270 **Figure 5: The DBN recapitulates the consequences of TAL1 and LYL1 single and double**
1271 **perturbations as seen *in vivo* and *in vitro*.** Computational prediction of gene expression
1272 patterns for the nine TFs of interest after perturbation of TAL1 **(a)**, LYL1 **(b)** or both **(c)**.
1273 Deletion of TAL1 or LYL1 on their own has no major consequences on the expression levels
1274 of the other eight TFs of the gene regulatory network, but simultaneous deletion of both TAL1
1275 and LYL1 caused changes in expression of several genes, mainly a decrease in *Gata2* and
1276 *Runx1*. This major disruption of the core GRN for blood stem/progenitor cells is therefore
1277 consistent with TAL1/LYL1 double knockout HSCs showing a much more severe phenotype
1278 than the respective single knock-outs. One thousand simulations were run for each perturbation
1279 to determine the TFs expression levels in a “cell population” by selecting expression levels at
1280 random time points after reaching its initial steady state. Expression levels of 0 resemble no
1281 expression, whereas expression levels of 1 stand for highest expression level that is possible in
1282 this system. The scale of the x-axis is linear. **(d)** Gene expression levels measured in single
1283 416b cells transfected with siRNA constructs against Tal1 or a control. The density plots of

1284 gene expression levels after perturbation of TAL1 indicate the relative number of cells (y-axes)
1285 at each expression level (x-axes). The scale of the x-axis is linear. The values indicate the
1286 results of the Wilcoxon rank-sum test: alterations to the expression profiles are indicated by the
1287 p-value (statistical significance: $p < 0.001$ for computational data and $p < 0.05$ for experimental
1288 data); substantial shifts in median expression level are indicated by the shift of median (SOM)
1289 (SOM >0.1 for computational data and >1 for experimental data). For details, see Fig. 5 –
1290 figure supplement 1; for full expression data, see figure 5 – source data.

1291

1292 **Fig. 5 - figure supplement 1: Significance tests for the computational and experimental**
1293 **data after TF perturbations.** To determine statistical significance the Wilcoxon rank-sum test
1294 was used. Alterations to the expression profiles are indicated by the p-value; with statistically
1295 significance defined as follows: $p < 0.001$ for computational data and $p < 0.05$ for experimental
1296 data. Significance of a substantial shift in median expression levels are as follows: shift of
1297 median >0.1 for computational data and >1 for experimental data (because of different scales).
1298 If the number for the shift of median is negative, the median of the perturbation data is smaller
1299 than that of the wild-type control; if the number is positive, the median of the perturbation is
1300 larger than that of the control. For simplicity, all significant changes are highlighted in red (p-
1301 value) and blue (shift of median).

1302

1303 **Figure 5 - figure supplement 2: Histogram plots showing the gene expression**
1304 **distributions of all nine genes of the network for the following perturbations: (a) LYL1**
1305 **down-regulation; (b) TAL1/SCL down-regulation; (c) LYL1 and TAL1/SCL down-regulation;**
1306 **(d) PU.1 down-regulation; (e) GFI1B up-regulation; and (f) AML-ETO9a simulation.**

1307

1308 **Figure 5 - source data: Raw and normalised data for the single cell gene expression**
1309 **experiments presented in this study:** 1) TAL1 down-regulation (related to Fig. 5 d), 2) PU.1
1310 down-regulation (related to Fig. 6 a), 3) GFI1B up-regulation (related to Fig. 6b) and 4) AML-
1311 ETO9a perturbation (related to Fig. 6 c)

1312

1313 **Figure 6: The DBN captures the transcriptional consequences of network perturbations.**

1314 **Left panel:** Computational prediction of gene expression after perturbation of specific TFs.
1315 1000 simulations were run for each perturbation to determine expression levels in a “cell
1316 population” (expression at 0 resembles no expression, whereas expression of 1 represents the
1317 highest possible expression level). The scale of the x-axis is linear. **Right panel:** Density plots
1318 of gene expression levels in single 416b cells after perturbation of specific TFs indicating the
1319 relative number of cells at each expression level. The scale of the x-axis is linear. The values
1320 indicate the results of the Wilcoxon rank-sum test: alterations to the expression profiles are
1321 indicated by the p-value (statistical significance: $p < 0.001$ for computational data and $p < 0.05$
1322 for experimental data); substantial shifts in median expression level are indicated by the shift of
1323 median (SOM) (SOM >0.1 for computational data and >1 for experimental data). For details,
1324 see Fig. 5 – figure supplement 1. **(a)** PU.1 down-regulation: (Left) Computational prediction of
1325 gene expression after PU.1 knockdown (Spi1 was set to 0 after reaching its initial steady state).
1326 (Right) Gene expression levels measured in single 416b cells transduced with shRNA
1327 constructs against shluc (wild-type) or shPU.1 (PU.1 knockdown). **(b)** GFI1B over-expression:
1328 (Left) Computational prediction of gene expression after over-expression of GFI1B (Gfi1b was
1329 set to 1 after reaching its initial steady state). (Right) Gene expression levels in single 416b
1330 cells transduced with a Gfi1b-expressing vector compared to an empty vector control (wild-
1331 type). **(c)** Consequences of the AML-ETO9a oncogene: (Left) Computational prediction of
1332 gene expression patterns after introducing the dominant-negative effect of the AML-ETO9a

1333 oncogene (Runx1 was fixed at the maximum value of 1 after reaching its initial steady state
1334 and in addition all Runt binding sites were set to have a repressive effect). (Right) Gene
1335 expression levels measured in single 416b cells transduced with an AML-ETO9a expressing
1336 vector fused to mCherry. mCherry positive cells were compared to mCherry negative cells
1337 (wild-type).

1338

1339 **Figure 6 – figure supplement 1: Summary of all computational simulations for**
1340 **perturbations of one or two TFs.** The results for a total of 162 simulations are shown. The
1341 data can be accessed using the embedded hyperlinks. The y-axes show the number of cells and
1342 the x-axes the relative expression level. Blue curves represent wild-type data and red curves
1343 represent perturbation data.

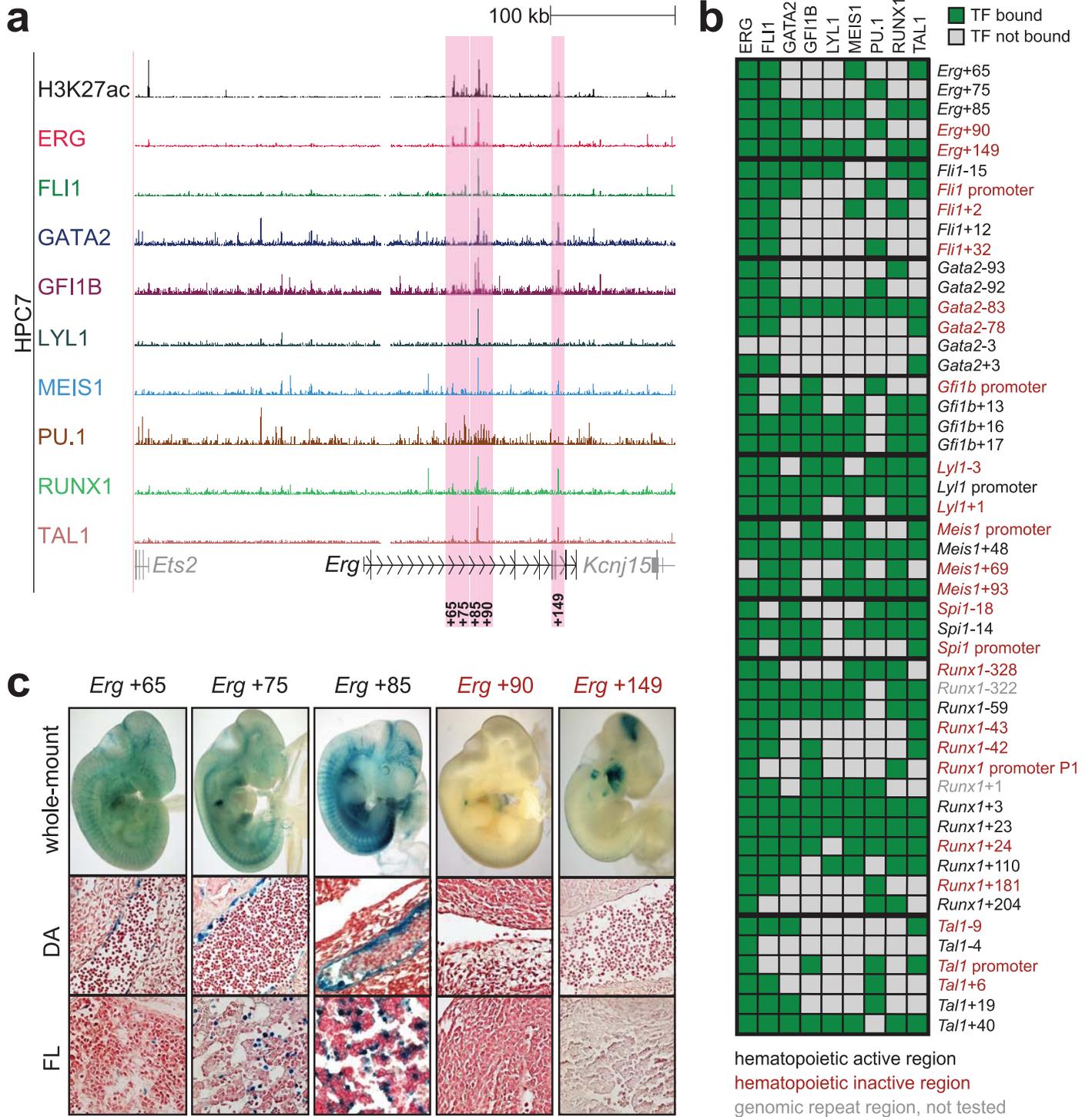


Figure 1

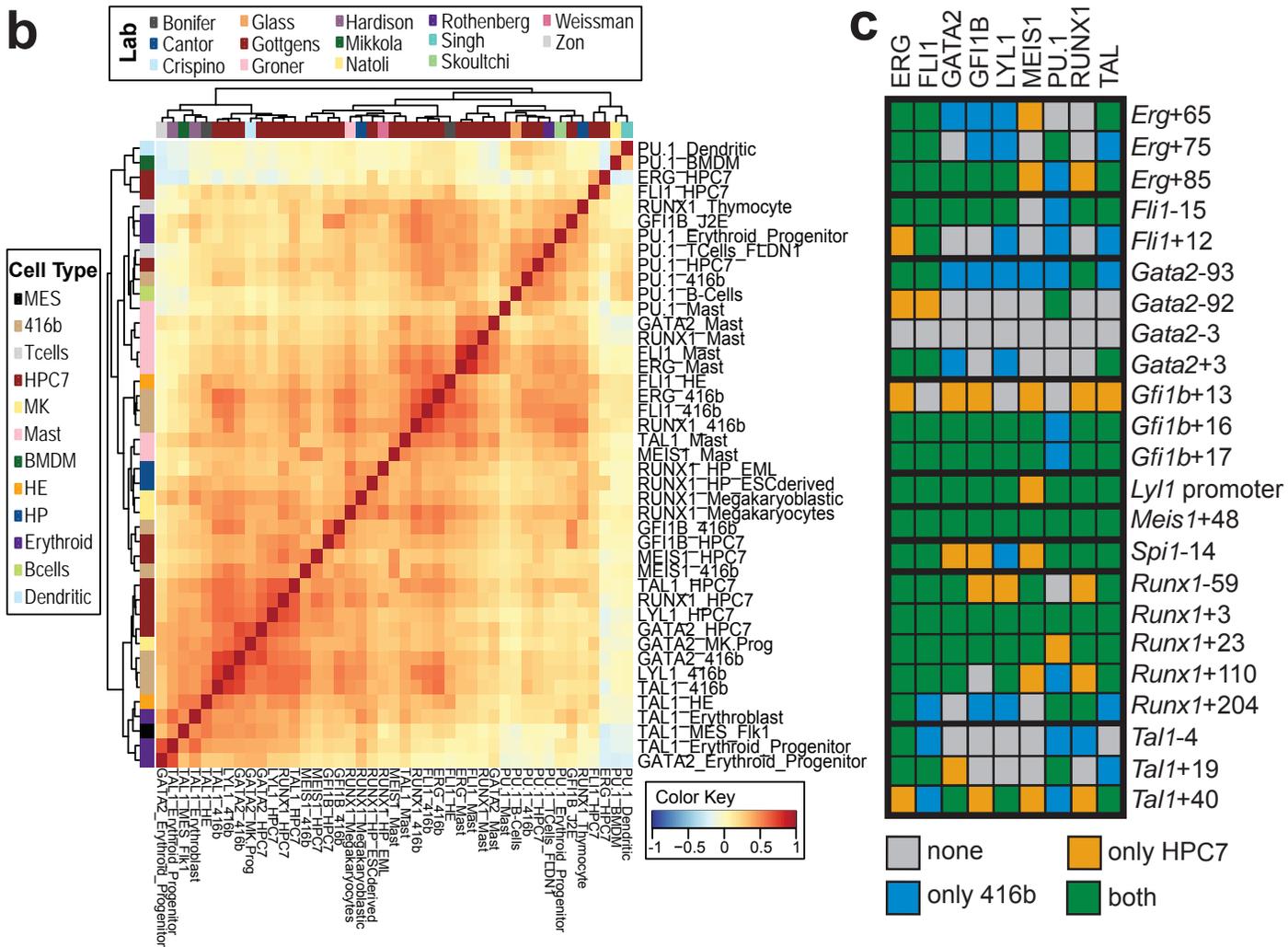
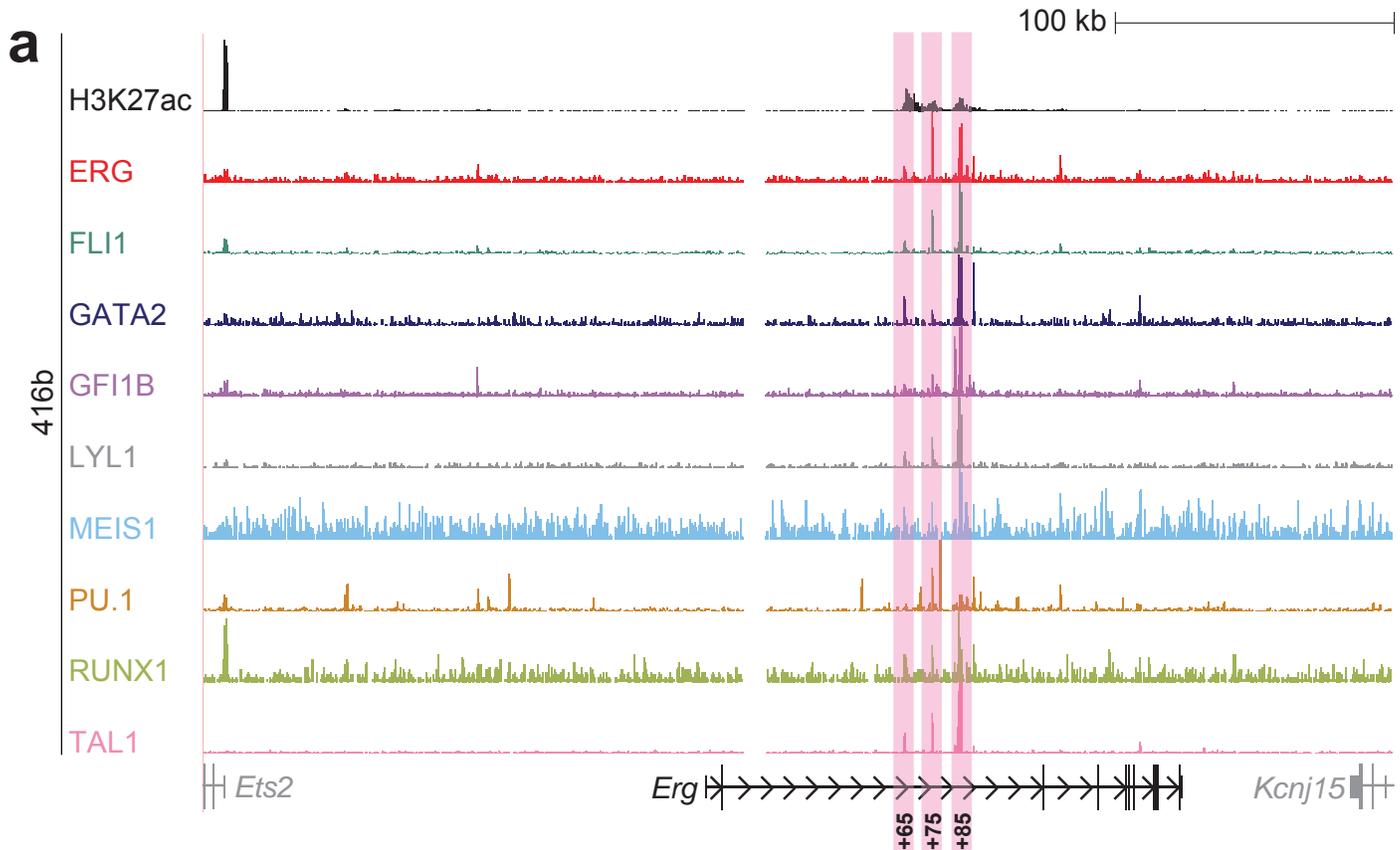


Figure 2

a *Erg+65*

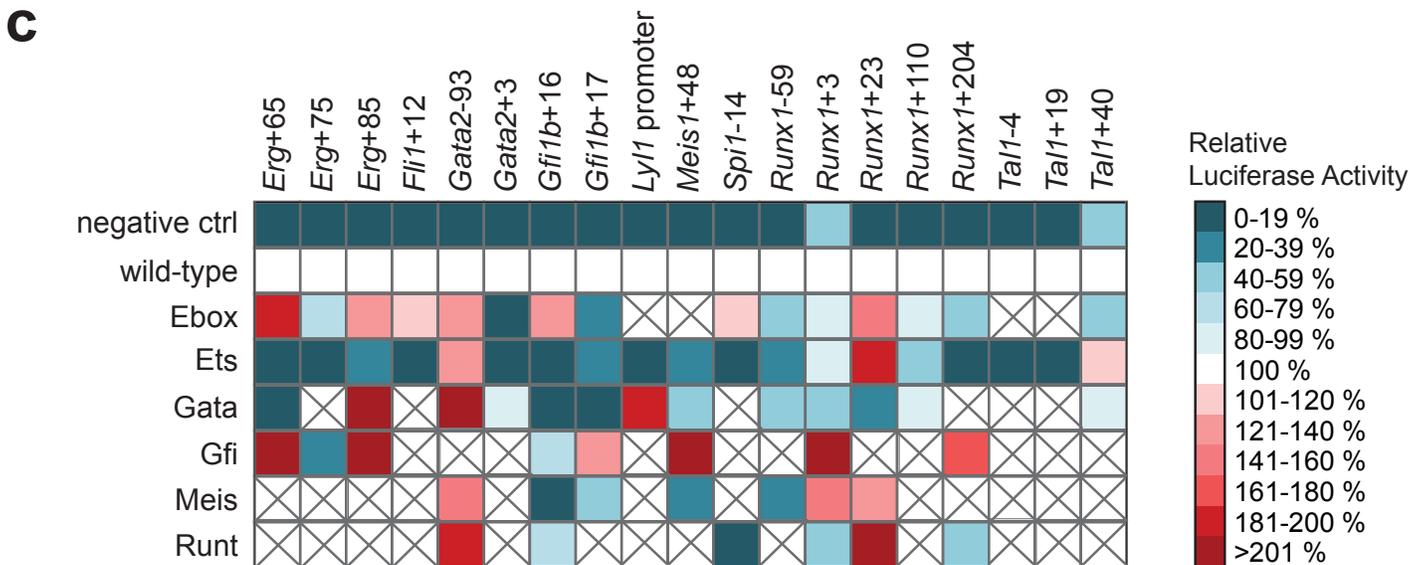
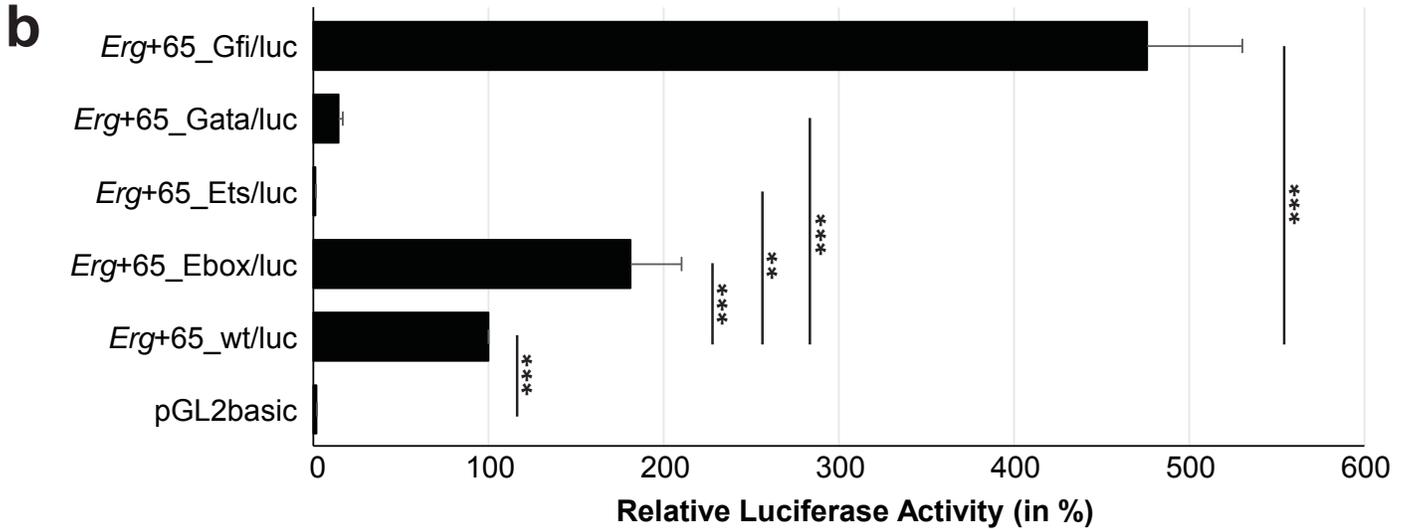
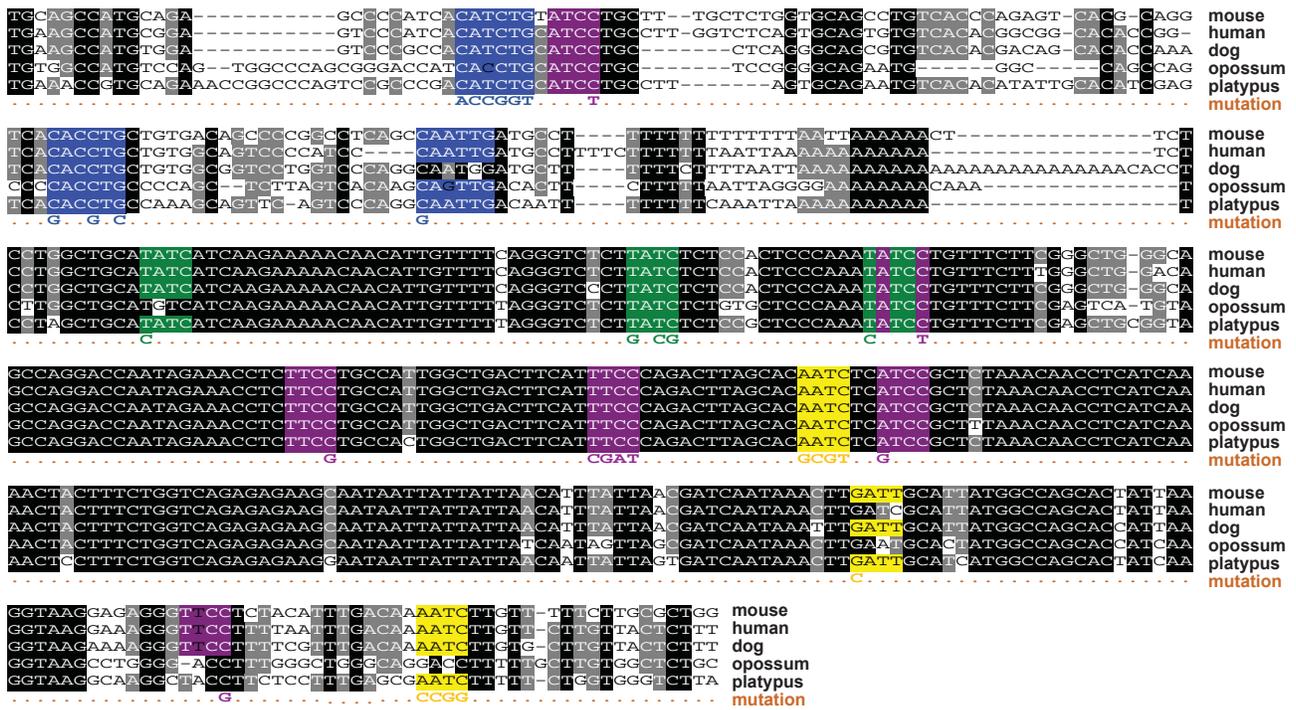


Figure 3

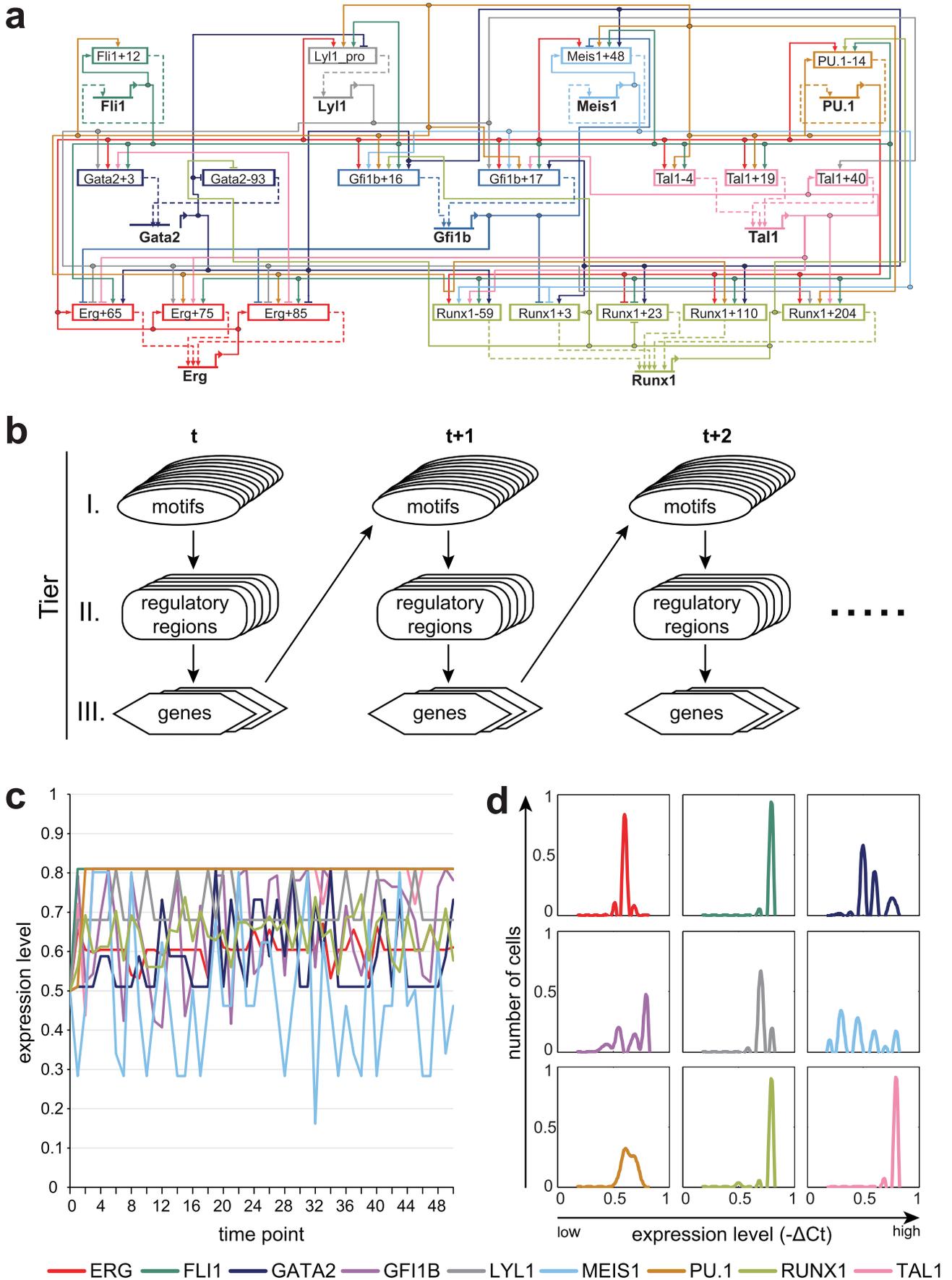
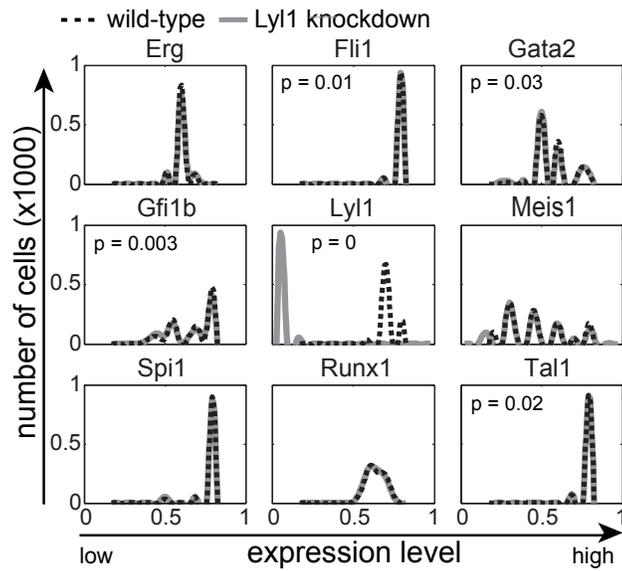
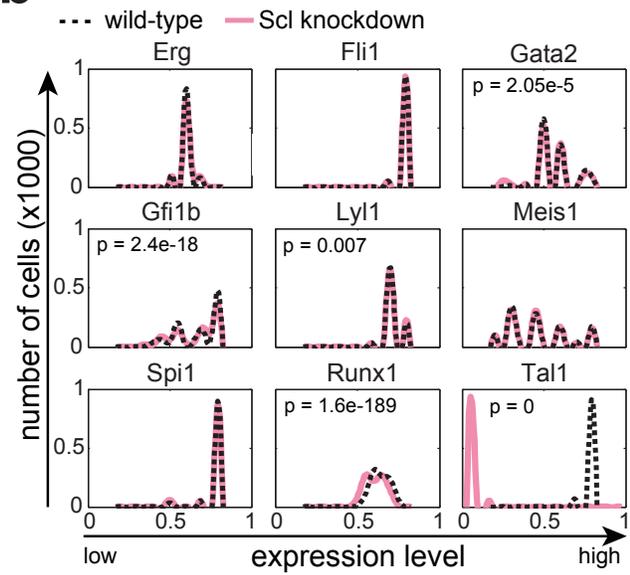
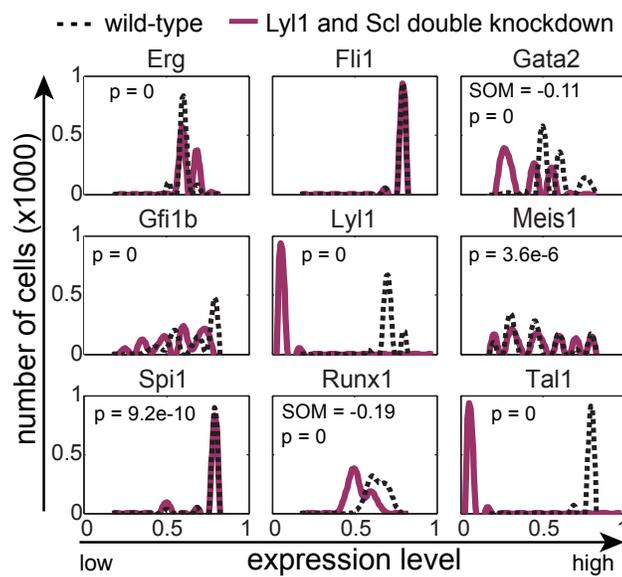
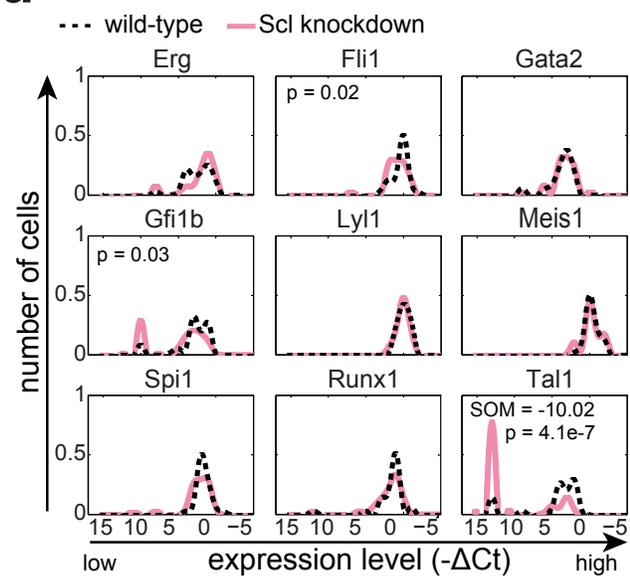


Figure 4

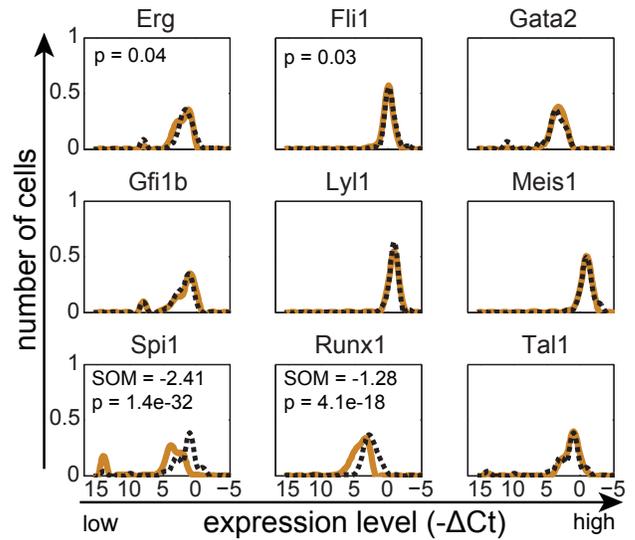
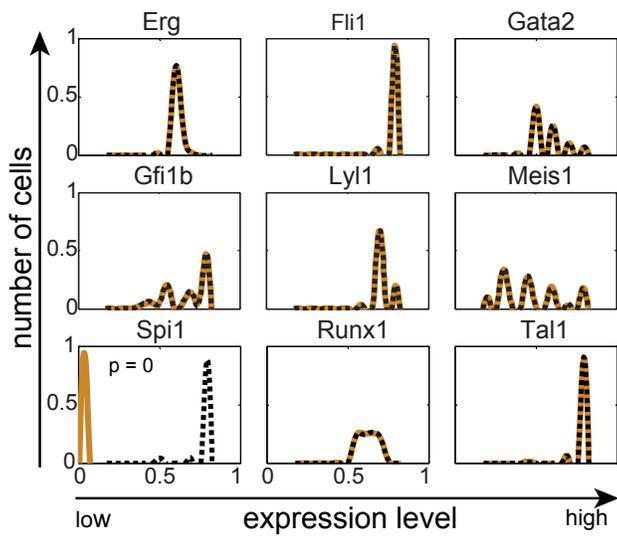
a LYL1 down**b TAL1 down****c LYL1 and TAL1 down****d TAL1 down (experimental)****Figure 5**

Computational

Experimental

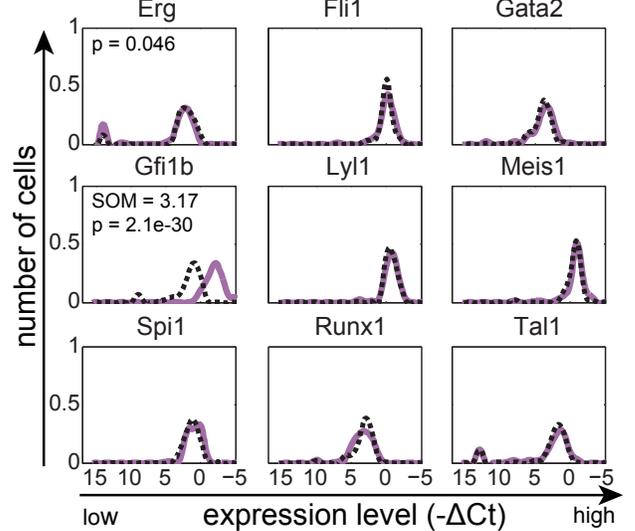
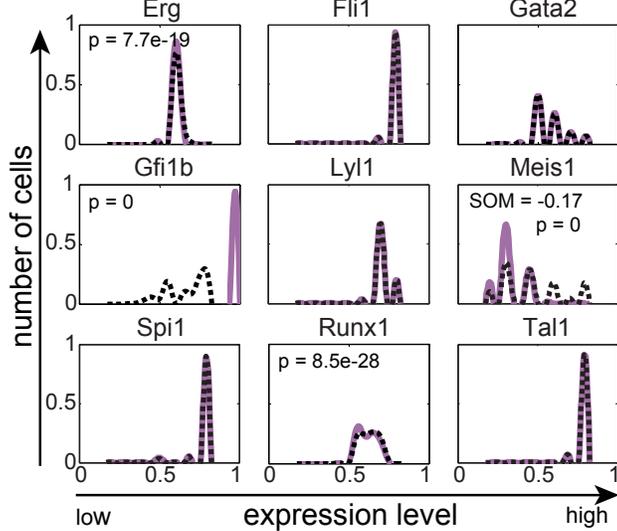
a PU.1 down

--- wild-type — PU.1 knockdown



b GFI1B up

--- wild-type — GFI1B over-expression



c AML1-ETO9a

--- wild-type — AML-ETO9a simulation/expression

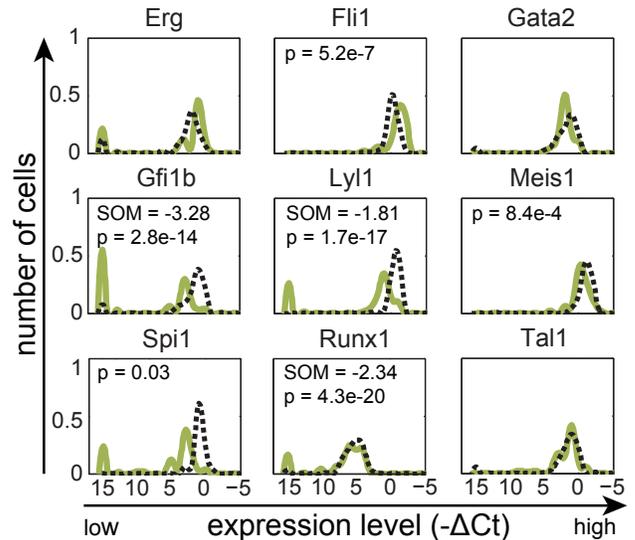
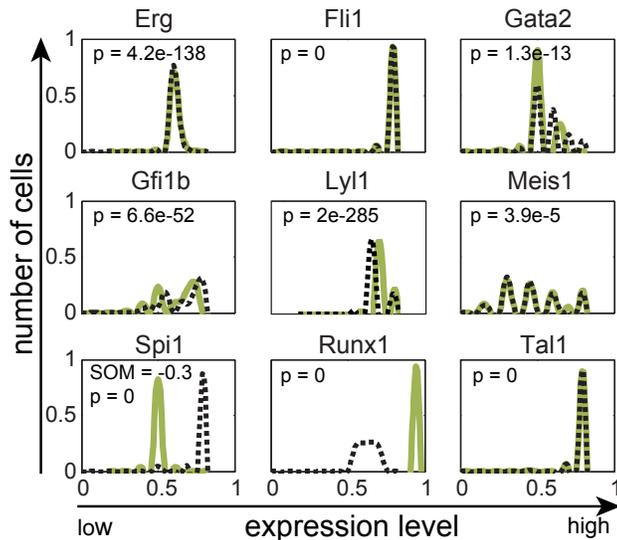


Figure 6