

# 1 **Coordinated Changes in Gene Expression Kinetics Underlie both**

## 2 **Mouse and Human Erythroid Maturation**

3 Melania Barile<sup>1,2</sup>, Ivan Imaz-Rosshandler<sup>1,2</sup>, Isabella Inzani<sup>3</sup>, Shila Ghazanfar<sup>4</sup>, Jennifer Nichols<sup>2,5</sup>, John  
4 C. Marioni<sup>4,6,7</sup>, Carolina Guibentif<sup>1,2,8,\*</sup>, Berthold Göttgens<sup>1,2,\*</sup>

5

- 6 1. Department of Haematology, University of Cambridge, CB2 0AW Cambridge, UK
- 7 2. Wellcome-Medical Research Council Cambridge Stem Cell Institute, University of Cambridge,  
8 CB2 0AW Cambridge, UK
- 9 3. University of Cambridge Metabolic Research Laboratories and MRC Metabolic Diseases Unit,  
10 CB2 0QQ Cambridge, UK
- 11 4. Cancer Research UK Cambridge Institute, University of Cambridge, CB2 0RE Cambridge, UK
- 12 5. Department of Physiology, Development and Neuroscience, University of Cambridge, CB2 3DY  
13 Cambridge, UK
- 14 6. Wellcome Sanger Institute, Wellcome Genome Campus, CB10 1SA Cambridge, UK
- 15 7. European Molecular Biology Laboratory, European Bioinformatics Institute (EMBL-EBI),  
16 Wellcome Genome Campus, CB10 1SD Cambridge, UK
- 17 8. Sahlgrenska Center for Cancer Research, Department of Microbiology and Immunology,  
18 University of Gothenburg, 413 90 Gothenburg, Sweden

19

20 \* Co-corresponding authors

21 BG: [bg200@cam.ac.uk](mailto:bg200@cam.ac.uk)

22 CG: [carolina.guibentif@gu.se](mailto:carolina.guibentif@gu.se)

23

24

25

## 26 **Abstract**

27 Single cell technologies are transforming biomedical research, including the recent demonstration  
28 that unspliced pre-mRNA present in single cell RNA-Seq permits prediction of future expression states.  
29 Here we applied this ‘RNA velocity concept’ to an extended timecourse dataset covering mouse  
30 gastrulation and early organogenesis. Intriguingly, RNA velocity correctly identified epiblast cells as  
31 the starting point, but several trajectory predictions at later stages were inconsistent with both real  
32 time ordering and existing knowledge. The most striking discrepancy concerned red blood cell  
33 maturation, with velocity-inferred trajectories opposing the true differentiation path. Investigating  
34 the underlying causes revealed a group of genes with a coordinated step-change in transcription, thus  
35 violating the assumptions behind current velocity analysis suites, which do not accommodate time-  
36 dependent changes in expression dynamics. Using scRNA-Seq analysis of chimeric mouse embryos  
37 lacking the major erythroid regulator *Gata1*, we show that genes with the step-changes in expression  
38 dynamics during erythroid differentiation fail to be up-regulated in the mutant cells, thus underscoring  
39 the coordination of modulating transcription rate along a differentiation trajectory. In addition to the  
40 expected block in erythroid maturation, the *Gata1* chimera dataset revealed induction of PU.1 and  
41 expansion of megakaryocyte progenitors. Finally, we show that erythropoiesis in human fetal liver is  
42 similarly characterized by a coordinated step-change in gene expression. By identifying a limitation of  
43 the current velocity framework coupled with *in vivo* analysis of mutant cells, we reveal a coordinated  
44 step-change in gene expression kinetics during erythropoiesis, with likely implications for many other  
45 differentiation processes.

46 (244 words)

47

## 48 **Introduction**

49 Cellular differentiation into diverse cell types underpins all metazoan development. Moreover, cellular  
50 differentiation processes are also crucial for stem cell-mediated tissue maintenance, and their  
51 perturbation has been implicated in ageing-associated regenerative failure as well as malignant  
52 transformation (Akunuru and Geiger, 2016; Schultz and Sinclair, 2016). Since cellular differentiation  
53 decisions are made at the level of individual cells, elucidation of the underlying molecular mechanisms  
54 requires the use of single cell approaches. It is no surprise therefore that recent innovations in single  
55 cell molecular profiling technologies have been embraced rapidly by developmental and stem cell  
56 biologists, with complete single cell gene expression maps now available for developing embryos of  
57 several model organisms (Gerber et al., 2018; Mahadevaiah et al., 2020; Wagner et al., 2018 reviewed

58 in Ton et al., 2020)), as well as large-scale datasets covering adult tissue homeostasis (Borrett et al.,  
59 2020; Dahlin et al., 2018; Weinreb et al., 2020).

60

61 Comprehensive molecular profiling necessarily entails the generation of snapshot data, because cells  
62 need to be fixed to examine their molecular content. This in turn represents a major drawback for the  
63 study of differentiation processes, which commonly occur over extended timeframes via complex  
64 trajectories underpinned by intricate decision-making processes. Much excitement was therefore  
65 generated by a recent seminal study (La Manno et al., 2018), which demonstrated that unspliced pre-  
66 mRNA present in scRNA-Seq datasets can be exploited to predict likely future expression states. This  
67 so-called RNA velocity concept is based on the notion that the ratio between unspliced and spliced  
68 RNA differs depending on whether a gene is in the process of being up- or downregulated. During  
69 upregulation, there is a relative increase in newly transcribed unspliced RNA, with the converse  
70 occurring during downregulation. The RNA velocity framework has rapidly gained traction across the  
71 wider single cell community, being applied across multiple experimental systems (Kanton et al., 2019;  
72 Zhang et al., 2019; Zhou et al., 2019), and also extended as part of the scVelo analysis suite (Bergen et  
73 al., 2020), which allows inclusion of genes whose transcript levels are not in steady state.

74 One system where the RNA velocity concept has particular potential is erythropoiesis, the process  
75 whereby oxygen-transporting red blood cells are generated from multipotent haematopoietic  
76 progenitors. Research into the transcriptional control processes of erythropoiesis led to several  
77 paradigmatic discoveries, including the dissection of distal transcriptional control elements (Grosveld  
78 et al., 1987; Higgs et al., 1990; Mettananda et al., 2016), as well as antagonistic transcription factor  
79 pairings as executors of lineage choice in multipotent progenitors (Zhang et al., 1999). During  
80 embryogenesis, a first so-called primitive wave of erythropoiesis occurs in the yolk sac, followed by a  
81 second definitive wave, initiated also in the yolk sac, then predominantly in the fetal liver and later in  
82 the adult bone marrow (McGrath and Palis, 2008). The zinc finger protein Gata1 represents the  
83 archetypal erythroid transcription factor, and is required for the maturation of both primitive and  
84 definitive erythroid cells (Fujiwara et al., 1996; Gutierrez et al., 2008; Pevny et al., 1995; Pevny et al.,  
85 1991), as well as megakaryocyte maturation (Shivdasani et al., 1997). However, the precise molecular  
86 processes affected by Gata1 deletion in early embryonic erythropoiesis have remained obscure,  
87 principally because conventional biochemical methods are unsuitable for the very small number of  
88 cells present at these early developmental stages.

89

90 Here, we have applied RNA velocity to a recently published scRNA-Seq dataset of nine sequential  
91 timepoints, spaced 6 hours apart, which encompass mouse gastrulation and early organogenesis  
92 (Pijuan-Sala et al., 2019). We observed that some of the inferred trajectories are incompatible with  
93 the existing biological knowledge, as well as the real time ordering derived from the sequential  
94 sampling timepoints. For erythroid differentiation in particular, we show that failure of the Velocity  
95 framework is due to a concerted increase in transcription rate of a subset of erythroid genes, midway  
96 through the red blood cell maturation trajectory. Analysis of *Gata1* chimeric embryos underscores  
97 the concerted nature of this expression boost, consistent with the notion that such concerted  
98 upregulation events may be a feature of stabilizing a given differentiated cellular state.

99

100

## 101 **Results**

### 102 **Limitations of RNA velocity trajectory inference at organismal scale**

103 To evaluate RNA velocity-based trajectory inference with a complex dataset, we applied the scVelo  
104 analysis pipeline (Bergen et al., 2020) to a recently reported timecourse scRNA-Seq dataset covering  
105 mouse gastrulation and early organogenesis. This mouse gastrulation atlas contains approximately  
106 120,000 single cell transcriptomes across nine sequential timepoints covering 37 major cell types  
107 (Pijuan-Sala et al., 2019). Prior to scVelo analysis, we removed extraembryonic ectoderm and  
108 extraembryonic endoderm cells, as they derive from early lineage branching events that are not  
109 covered in this dataset. We first applied scVelo to the normalised and batch corrected count matrix  
110 across all embryonic stages (Figure 1A). We observed that scVelo correctly identifies the epiblast  
111 population as the origin of the global differentiation processes that occur during gastrulation and early  
112 organogenesis. In relation to the more differentiated cell types however, there were several instances  
113 where scVelo had difficulty in capturing some of the highly complex differentiation events that occur  
114 across the entire embryo. For instance, scVelo predicted that E8.0 allantois and mesenchyme cell-  
115 types give rise to mesodermal cells from earlier timepoints rather than the E8.25/E8.5 allantoic and  
116 mesenchymal cells. Another inconsistency occurred with E8.0-E8.25 endoderm cells, which were  
117 predicted to give rise to E6.5-E7 visceral endoderm, rather than the other way round. Most  
118 noteworthy, scVelo failed to recapitulate the erythropoiesis branch, where it predicts a backwards  
119 differentiation from later to earlier populations. We next repeated this analysis using data from each  
120 individual time-point (Figure 1B; shown are E7.5 and E8.5). We saw that the pipeline accurately

121 recapitulates known biological trajectories up to E7.5, but observed the same inconsistency from  
122 E7.75 to E8.5, with scVelo arrows pointing backwards.

123 Taken together therefore, we have identified that for erythroid development, the output of scVelo is  
124 inconsistent with the timecourse information gathered from the experimental design of the  
125 gastrulation atlas.

126

### 127 **Unspliced sequence reads help to discriminate between cell types**

128 We next asked whether this issue is due to a general lack of biologically meaningful information  
129 captured in the unspliced reads.

130 To this end, we exploited two variance-based dimensionality reduction methods, Principal Component  
131 Analysis (PCA) and Multi-Omics Factor Analysis (MOFA; Argelaguet et al., 2020), to interrogate how  
132 much inter-population variability is explained by the spliced and unspliced information layers, whether  
133 considered separately or together. Upon comparing PC1 and PC2 (or MOFA Factors 1 and 2), in  
134 addition to the expected lineage separation obtained using the spliced reads (Figure 2A, left panel),  
135 we could also observe a degree of lineage separation when using the unspliced reads alone (Figure  
136 2A, middle panel). In addition, we saw a qualitatively improved separation of the different lineages  
137 when spliced and unspliced information is used in combination (Figure 2A, right panel; see  
138 Supplementary Figure 1 for further components/factors). Moreover, the MOFA factors account for  
139 16% of variation in the spliced data and 4% of the of variation in unspliced data (Figure 2Bi).  
140 Interestingly, a closer look at the MOFA pre-processing and final outcome showed a minor overlap of  
141 genes that are highly variable with respect to spliced or unspliced counts (Figure 2Bii) and a different  
142 weight contributed by the two layers to the final factors (Figure 2Biii).

143 Multiomics factor analysis therefore not only demonstrates that the unspliced reads in the  
144 gastrulation atlas dataset contain biologically relevant information, but also suggests that integrated  
145 analysis of spliced and unspliced reads may more broadly facilitate the interpretation of complex  
146 scRNA-Seq datasets.

147

### 148 **Analysis of unspliced reads reveals complex expression kinetics**

149 Having confirmed the utility of unspliced reads, we next explored whether the inability to recover real-  
150 time progression in whole embryo trajectory inference using scVelo might be related to the  
151 assumptions made by the current RNA velocity analysis tools. The derivation of gene-specific

152 expression kinetics underpins the scVelo analysis pipeline, as illustrated by so-called phase plots that  
153 depict the amounts of spliced versus unspliced reads within a population of cells (Bergen et al., 2020).  
154 If a gene is upregulated during a differentiation timecourse, cells will be placed above the diagonal  
155 between no expression and maximum expression due to the relatively larger amount of newly  
156 produced pre-mRNA during the gene induction process, while the converse is true for downregulated  
157 genes (Figure 3A). Both of these scenarios are readily captured by scVelo, with the predicted vectors  
158 of differentiation agreeing with the actual temporal progression. If a given gene however experiences  
159 an increase in transcription rate midway through a differentiation timecourse, the sudden increase in  
160 unspliced pre-mRNA will result in a phase plot that may be wrongly classified by scVelo, with predicted  
161 vectors of differentiation diametrically opposed to the true direction of differentiation (Figure 3A).  
162 This is indeed what we observed when inspecting the phase plots of the scVelo driver genes (top-  
163 likelihood genes, Supplementary Table 1), which display a steep increase of unspliced counts in the  
164 Erythroid 3 population, leading to a reverse velocity prediction, progressing from Erythroid 3 to earlier  
165 populations (Supplementary Figure 2A).

166 We next set out to identify all genes exhibiting this rapid increase in expression levels in the Erythroid  
167 3 population (Figure 3B). After fitting a linear regression through each population and each gene and  
168 testing whether the inferred slopes reflected the expected order based on biological knowledge, we  
169 found 73 such genes, which we termed Multiple Rate Kinetics or MURK genes. These genes included  
170 *Smim1* and *Hba-x*, where we could confirm an increase in expression kinetics using phase plots (Figure  
171 3C).

172 Having identified a set of genes with a coordinated increase in expression rate midway through  
173 erythropoiesis, we next asked what function these genes might play in the broader transcriptional  
174 program of red blood cell maturation. Visual inspection of the gene list revealed it to contain  
175 archetypal red blood cell genes including the globin genes *Hba-x*, *Hbb-a1*, *Hba-a2*, *Hbb-bt*, *Hbb-bh1*,  
176 *Hbb-y* (Supplementary table 2). Unsupervised gene ontology analysis confirmed that biological  
177 functions essential for red blood cells were highly enriched, including “gas transport” and “heme  
178 biosynthetic process” (Figure 3D, Supplementary Figure 2B).

179 We next removed this set of MURK genes and recalculated the RNA velocity inferred trajectories. As  
180 can be seen in Figure 3E, inferred vectors of differentiation are now in good agreement with the real  
181 time progression of erythropoiesis, with the only discrepancies occurring towards the end,  
182 presumably due to the fact there are no future expression states to project into beyond E8.5, the last  
183 timepoint sampled in the analysed dataset.

184 The scVelo suite also calculates a so-called latent time, which represents the pseudotime ordering  
185 hidden in the spliced and unspliced dynamics, and is more powerful than previously described  
186 pseudotime inferring approaches since it incorporates both the gene dynamics and the spliced and  
187 unspliced information (Bergen et al., 2020). Using the full gene set, the latent time calculation for the  
188 erythroid lineage is contrary to the know progression of erythroid differentiation (Figure 3E left panels,  
189 Supplementary Figure 2C, left panels). By contrast, removing the MURK genes results in a latent time  
190 prediction that is not only consistent with the major axis of erythropoiesis, but also identifies the two  
191 sequential inputs described previously (Pijuan-Sala et al., 2019), namely an early wave directly from  
192 posterior mesoderm as well as a second wave coming from yolk sac hemogenic endothelium (see  
193 Figure 3E, Supplementary Figure 2C, right panels).

194 Taken together therefore, this analysis shows that inconsistent RNA velocity-inferred trajectories can  
195 be remedied by the removal of genes with complex expression kinetics.

196

#### 197 **Erythroid Multiple Rate Kinetics genes are essential for red blood cell function**

198 To corroborate upregulation of our identified MURK genes during erythropoiesis, we interrogated a  
199 previously published dataset with transcriptomic analysis of a loss of function model for the  
200 erythropoiesis master-regulator *Gata1* (Wu et al., 2011). *In vitro* differentiation of *Gata1* knock-out  
201 embryonic stem cells over-expressing human *BCL2* can produce permanently self-renewing immature  
202 erythroid progenitor cell lines. One such model, G1ER, contains a tamoxifen-inducible *Gata1*  
203 transgene, the activation of which triggers erythroid maturation (Tsang et al., 1997; Weiss et al., 1997;  
204 Figure 4A). Microarray-based differential gene expression was performed, comparing the uninduced  
205 and induced conditions (Wu et al., 2011). 70 of our 73 MURK genes overlapped with the genes  
206 identified by this microarray-based comparison. Of those, 52 were upregulated, of which 45 showed  
207 strong upregulation, 4 were downregulated, and 4 showed no change in expression following  
208 induction of *Gata1* in the G1ER system, demonstrating a highly significant overlap of our identified  
209 MURK genes with the G1ER-induced genes ( $p < 10^{-11}$ ; see Figure 4B).

210 Our newly identified erythropoietic MURK genes therefore perform key roles in red blood cell  
211 function, and their upregulation was validated in an independent model of red blood cell maturation.

212

#### 213 **scRNA-Seq of mouse chimeras reveals the early cellular defects *Gata1* loss of function**

214 The G1ER cell line represents an *in vitro* model, and the published differential gene expression data  
215 were from bulk microarray profiling, thus precluding any analysis of single-cell gene expression  
216 kinetics. We therefore turned to our recently reported Chimaera-Seq approach, whereby scRNA-Seq  
217 is coupled with mouse chimeric embryo technology, to define both cellular and molecular  
218 consequences of gene knock-outs *in vivo* (Pijuan-Sala et al., 2019). We used our standard embryonic  
219 stem cells (ESCs) expressing a constitutive tdTomato (tdTom) fluorescent marker gene to generate a  
220 *Gata1* knock-out line (see Methods). *Gata1*<sup>-</sup> tdTom<sup>+</sup> cells were injected into tdTom<sup>-</sup> wild-type  
221 blastocyst and transferred into pseudo-pregnant females, resulting in chimeric embryos that we  
222 harvested at E8.5. Six chimeric embryos were pooled, dissociated into a single-cell suspension, and  
223 tdTom<sup>+</sup> and tdTom<sup>-</sup> cell fractions were sorted for scRNA sequencing. We obtained 8420 tdTom<sup>-</sup> and  
224 7944 tdTom<sup>+</sup> cells passing quality control and assigned to a cell type, with an average of 4354 genes  
225 being detected per cell.

226 We then concatenated the chimera data with the Pijuan-Sala et al. (2019) reference dataset and  
227 mapped nearest neighbors (see Methods). We observed an overall homogeneous distribution of both  
228 mutant and wild-type fractions throughout the later time-points of the landscape, except for the  
229 erythroid branch. Indeed, we observed a block in the erythroid lineage of the mutant cells, which were  
230 over-represented in the start of the erythroid differentiation branch, while their wild-type  
231 counterparts were present throughout erythroid differentiation (Supplementary Figure 3).  
232 Identification of the nearest neighbours of chimeric cells within the reference dataset allowed their  
233 quick cell-type annotation, which we used to quantify the differences in the hemato-endothelial cell-  
234 type representation within the chimera fractions. This analysis confirmed a severe erythroid  
235 differentiation defect of the mutant cells (Figure 4C-E). When examining the reference dataset  
236 sampled-time point of the chimera nearest neighbours we also observed a temporal shift within the  
237 erythroid lineage, with tdTom<sup>+</sup> mutant cells mapping to earlier time-points than their wildtype tdTom<sup>-</sup>  
238 counterparts, further confirming a developmental block of the mutant cells (Figure 4D, E). In addition,  
239 we observed that this erythroid defect was coupled with an over-representation of cells with a  
240 megakaryocyte signature (Figure 4C).

241 The newly generated *Gata1*<sup>-</sup> Chimaera-Seq data therefore not only recapitulated the expected block  
242 in erythroid maturation, but also revealed an expansion of the megakaryocytic lineage in the E8.5 yolk  
243 sac.

244

245 **The molecular program affected by *Gata1* loss in early embryos**

246 Although the role of Gata1 is well documented in developmental erythropoiesis (Fujiwara et al., 1996;  
247 Pevny et al., 1995), the early molecular defects of Gata1 loss of function *in vivo* had never been  
248 reported. The Gata1 Chimaera-Seq dataset therefore presented an opportunity to dissect the early  
249 molecular program controlled by Gata1 *in vivo*. Having registered a defect in erythroid differentiation  
250 and an increase in the megakaryocytic lineage population, we performed differential gene expression  
251 testing between the chimera mutant and wild-type cells in these clusters (Supplementary Table 3).

252 Regarding the megakaryocytic subset, we observed upregulation of progenitor markers *Kit*, *Gata2* and  
253 *Myb* in the *Gata1*<sup>-</sup> cells as well as lower expression of maturation genes for the megakaryocyte lineage  
254 *Gp5*, *Pf4*, *Mpl* and *Plek* (Figure 5A). Hyper-proliferative megakaryocyte progenitors, detected  
255 previously in *Gata1*<sup>-</sup> E12.5 fetal livers, led to compromised platelet function, and were suggested to  
256 originate in the yolk sac (Vyas et al., 1999). Our results showing over-production of megakaryocytic  
257 cells with impaired maturation characteristics in E8.5 *Gata1*<sup>-</sup> chimera yolk sacs support this notion,  
258 and importantly place the megakaryocytic defect within the very early phase of megakaryocyte  
259 formation.

260 Interestingly, all hemato-endothelial cell subsets displayed up-regulation of *Spi1* (coding for the PU.1  
261 transcription factor) in the *Gata1*<sup>-</sup> cell fraction compared to wild-type counterpart (FDR < 0.01; Figure  
262 5A). Given the previously reported Gata1-PU.1 cross-repression in adult bone marrow (Zhang et al.,  
263 1999) and in zebrafish embryonic hematopoiesis (Monteiro et al., 2011), we systematically assessed  
264 the effect of *Gata1* knockout in the mouse chimera lineages and observed that in *Gata1*<sup>-</sup> cells, *Spi1*  
265 was specifically up-regulated in all hematopoietic sub-clusters, with a stronger effect on Mk and Ery1  
266 subsets. (Supplementary Figure 3).

267 In the early erythroid subset, Ery1, we again noted that the mutant cells displayed increased  
268 expression of genes characteristic of a progenitor signature. Conversely, erythroid maturation  
269 hallmark genes such as *Hbb-bs* and *Gypa* were downregulated, along with the erythroid Gata1 target  
270 *Mllt3* (Pina et al., 2008; Figure 5A). GO-term enrichment analysis of genes downregulated in *Gata1*<sup>-</sup>  
271 Ery1 cells revealed biological processes essential to red blood cell function (Figure 5B). Furthermore,  
272 we also observed that 37% of the MURK genes identified in Figure 3 overlapped with these genes that  
273 fail to up-regulate in *Gata1*<sup>-</sup> erythroid cells (Figure 5C;  $p < 10^{-13}$ ).

274 In addition to the failure of inducing genes associated with erythroid maturation, single cell resolution  
275 molecular analysis also revealed a striking failure to downregulate genes associated with alternative  
276 lineage programs such as Pu.1, consistent with the notion that the earliest wave of primitive  
277 hematopoiesis produces erythroid cells, megakaryocytes and macrophages, with evidence for at least  
278 bipotential progenitor cells (Palis, 2016).

279

## 280 **The late erythroid increase in expression rate is downstream of Gata1 function**

281 Having generated the Chimaera-Seq single cell data for both wildtype and Gata1 knock-out cells, we  
282 next used the ratio of spliced/unspliced reads to explore differences in expression kinetics between  
283 the wildtype and mutant cells. As can be seen in Figure 5D, the previously defined MURK genes failed  
284 to display the increased rate of expression characteristic for the later stages of erythropoiesis in the  
285 mutant cells. The examples shown include the embryonic globin gene *Hbb-y*, as well as the *Smim1*  
286 gene coding for the Vel Blood Group Antigen (Storry et al., 2013) and the *Fam210b* gene, coding for a  
287 putative mitochondrial protein recently implicated in erythroid differentiation (Kondo et al., 2016;  
288 Figure 5D). This result confirms that the erythroid boost in expression forms part of the transcriptional  
289 program downstream of Gata1 function, although it does not demonstrate a direct regulatory role for  
290 Gata1. However, preliminary modelling analysis suggests that the change observed in MURK gene  
291 dynamics is due to altered transcription rates (see Supplementary Note), indicating a close association  
292 of the coordinated late erythroid increase in transcription rate with the molecular program  
293 downstream of Gata1.

294

## 295 **A coordinated increase of expression rate during human fetal liver erythropoiesis**

296 Having identified a coordinated increase in transcription rate during mouse yolk sac erythropoiesis,  
297 we next wanted to ascertain whether the same phenomenon could also be seen in human cells.  
298 Moreover, we were keen to explore an scRNA-Seq dataset generated by a different laboratory, to  
299 exclude any potential technical bias caused by our own experimental protocols. We therefore turned  
300 to a recently published comprehensive dataset of human fetal liver erythropoiesis (Popescu et al.,  
301 2019), and extracted the 49,388 cells annotated to the four clusters encompassing human fetal liver  
302 erythropoiesis. When calculating scVelo-based differentiation vectors as well as latent time using the  
303 full gene set (see methods), both were reversed (Figure 6A, left plots), consistent with the mouse yolk  
304 sac results. We therefore again ran our pipeline to discover genes with a potential increase in  
305 expression rate along the differentiation pathway. The resulting 377 genes again contained archetypal  
306 erythroid genes such as the hemoglobins (Figure 6B), with overall gene ontologies demonstrating a  
307 functional role in erythropoiesis (see Figure 6C); with in addition a high prevalence of cell-cycle related  
308 terms. We then recalculated both the scVelo differentiation vectors as well as latent time after  
309 removing the fetal liver MURK genes. This revealed scVelo vectors that were consistent with the  
310 expected developmental progression (see Figure 6A, right plots). This analysis therefore demonstrates

311 that complex expression kinetics apply broadly to erythropoiesis, and their identification can be used  
312 to amend the RNA velocity framework to prevent erroneous predictions.

313

## 314 **Discussion**

315 There is no doubt that single cell molecular profiling constitutes a transformative technology. It suffers  
316 however from the major drawback that cells need to be fixed in order to profile them, with the  
317 consequence that measurements are by necessity static snapshots. To decipher complex biological  
318 processes, however, temporal information is commonly required. The single cell RNA velocity concept  
319 raised the prospect of overcoming some of the limitations associated with static measurements, by  
320 providing a strategy that can infer future cellular states. The RNA velocity framework is based on an  
321 explicit model of transcriptional processes (transcription, splicing, degradation). The notion that  
322 physical parameters of gene expression can be deduced from single cell gene expression data had  
323 been explored before the single cell RNA velocity concept was introduced (Ezer et al., 2016; Kim and  
324 Marioni, 2013). However, the scVelo implementation provided an attractive framework for estimating  
325 gene-specific expression parameters by taking advantage of the spliced versus unspliced read counts  
326 across large cell populations (Bergen et al., 2020). Using erythropoiesis as an example, we show here  
327 that this current framework needs to be adapted to accommodate more complex expression kinetics.  
328 Importantly, our analysis revealed that sets of genes can show a coordinated increase in transcription  
329 rate along a differentiation pathway. Moreover, deletion of the key erythroid regulator Gata1  
330 abrogated this coordinated change in expression dynamics, thus revealing this increase in  
331 transcription rate as an important feature of erythropoiesis. It remains to be seen how dynamic  
332 changes in expression kinetics may play a role in other systems, where splicing and degradation as  
333 well as transcription rates may also be time-dependent.

334 Application of the single cell RNA velocity concept has commonly been “confirmatory”, whereby a  
335 differentiation path proposed by other means was shown to be consistent with RNA velocity  
336 inference. When we applied the RNA velocity framework to the entire mouse gastrulation atlas, some  
337 inferred vectors of differentiation agreed with our current understanding of developmental biology,  
338 but others disagreed. Deeper interrogation of predictions that conflicted with our current  
339 understanding of erythropoiesis showed that the RNA velocity predictions could not be correct, not  
340 only because they ran counter to the known expression changes that accompany red blood cell  
341 differentiation, but also because they contradicted the real-time sampling of the data. Our results thus  
342 highlight certain limitations of the current implementation of this framework for identification of  
343 novel trajectories. Importantly however, it is through our observation of the inconsistent predictions

344 that we were led to identify the previously unrecognized dynamic nature of the transcriptional control  
345 of erythropoiesis. Moreover, it is plausible that coordinated increases in transcription rate midway  
346 through a differentiation process may operate more widely, as a powerful mechanism for stabilising  
347 a cell state. Our extension to the scVelo implementation reveals the presence of such time-dependent  
348 changes of gene expression parameters and retrieves the concerned MURK genes in developmental  
349 trajectories of interest.

350 As to the precise mechanisms, at this stage we can only confidently assert that this process occurs  
351 downstream of Gata1 during erythropoiesis. Of note, comprehensive analysis of the G1ER erythroid  
352 differentiation model has shown that Gata1-induced maturation triggers increased  
353 enhancer/promoter interactions for upregulated genes, and that the most highly enriched motif in  
354 the promoters of these genes are GATA sites (Liu et al., 2020). These observations are therefore  
355 consistent with the lineage-determining function of Gata1 involving a coordinated increase in  
356 expression kinetics of a set of genes important for red blood cell function.

357 Our observations regarding the Gata1 knock-out phenotype also warrant some discussion. With  
358 embryonically lethal phenotypes such as Gata1 knock-out, conventional analysis tends to be  
359 somewhat limited, since the embryos are dead because they have no red blood cells. By contrast, the  
360 Chimaera-Seq assay enables both quantification of cell numbers as well as characterisation of their  
361 molecular profiles. Moreover, there are no secondary effects caused by the dying embryo, because  
362 the wildtype host cells rescue overall fetal development, thus allowing a focussed analysis of cell-  
363 intrinsic molecular defects. One noteworthy observation from our data is that erythroid  
364 differentiation proceeds substantially beyond the stage where *Gata1* expression itself is first initiated,  
365 but fails to proceed to the late erythroid phase where expression of canonical red blood cell genes is  
366 greatly upregulated. However, gene expression prior to the differentiation block is not normal. In  
367 particular, we observed increased *Spi1/Pu.1* in the Gata1 knock-out cells, consistent with the  
368 previously reported (Zhang et al., 1999) but also disputed (Hoppe et al., 2016) antagonistic  
369 relationship between Gata1 and Pu.1.

370 Within haematopoiesis, Pu.1 is recognised as a key regulator of myeloid and T-cell lineages, but not  
371 erythroid cells, even though a role in the proliferation of immature erythroid progenitors has been  
372 reported (Choe et al., 2010; reviewed in Carotta et al., 2010). Upregulation of Pu.1 in our immature  
373 *Gata1* knock-out cells therefore suggests that these cells of the primitive haematopoietic lineage  
374 represent progenitors with multilineage potential, rather than being restricted to just the red cell  
375 lineage. Further evidence for this notion is provided by our observation that the reduction in erythroid  
376 cells in the *Gata1* knock-out is accompanied by an increase in megakaryocyte progenitors, consistent

377 with a model whereby Gata1 levels influence the lineage choice decisions of a multipotent progenitor  
378 cell. Live cell tracking studies have suggested that the primary role of Gata1 and Pu.1 may be fate  
379 stabilization rather than fate choice (Hoppe et al., 2016). The increase in transcription rate of erythroid  
380 genes downstream of Gata1 would cohere with stabilizing the erythroid fate, thus suggesting that our  
381 results are consistent with roles in both fate choice and fate stabilization.

382 Our observation of an expanded pool of megakaryocyte progenitors may also be of direct relevance  
383 to our understanding of the pre-leukaemic transient myeloproliferative disease (TMD) that is  
384 prevalent in newborns with trisomy 21 (Roberts et al., 2013). TMD is thought to arise when a fetal  
385 specific haematopoietic progenitor cell with trisomy 21 acquires a partial loss of function mutation in  
386 *GATA1*, resulting in a short form of GATA1 (GATA1s). TMD is characterized by expansion of immature  
387 megakaryocyte progenitors, and in 10 to 20% of cases transforms into malignant acute  
388 megakaryoblastic leukaemia (reviewed in Bhatnagar et al., 2016)). Over-expression of GATA1s in  
389 mouse models resulted in the identification of mid-gestation fetal liver megakaryocyte progenitors as  
390 uniquely sensitive to this mutant GATA1s form compared to their adult bone marrow counterparts (Li  
391 et al., 2005). The over-represented population of immature megakaryocytic progenitors in our E8.5  
392 *Gata1*<sup>-</sup> chimeras may correspond to the developmental emergence of this transient precursor, TMD-  
393 initiating cell, in the yolk sac.

394 Taken together, this study reports how the RNA velocity framework can be extended to delve into the  
395 transcriptional mechanisms of tissue differentiation, complemented with single cell resolution and *in*  
396 *vivo* analysis of Gata1 function, which revealed a number of previously unknown facets of this  
397 canonical regulator of red blood cell development.

398

### 399 **Acknowledgements**

400 We would like to thank Prof. Fabian Theis and Volker Bergen for discussions and valuable input on the  
401 scVelo implementation. We thank William Mansfield and the Gurdon Institute animal facility for  
402 blastocyst injections, the Flow Cytometry Core Facility at CIMR for cell sorting, Katarzyna Kania and  
403 the CRUK-CI genomics core for preparing the 10X libraries and for sequencing. Research in the authors'  
404 laboratories is supported by the Wellcome Trust, MRC, CRUK, Blood Cancer UK, NIH-NIDDK, the  
405 Sanger-EBI Single Cell Centre; by core support grants by the Wellcome Trust to the Cambridge Institute  
406 for Medical Research and Wellcome Trust-MRC Cambridge Stem Cell Institute; and by core funding  
407 from Cancer Research UK and the European Molecular Biology Laboratory. C.G. was funded by the  
408 Swedish Research Council (2017-06278), I.I. was funded by a British Heart Foundation studentship

409 (FS/18/56/35177), S.G. was supported by a Royal Society Newton International Fellowship  
410 (NIF\R1\181950). This work was funded as part of a Wellcome Strategic Award (105031/D/14/Z)  
411 awarded to Wolf Reik, Berthold Göttgens, John Marioni, Jennifer Nichols, Ludovic Vallier, Shankar  
412 Srinivas, Benjamin Simons, Sarah Teichmann, and Thierry Voet.

413

#### 414 **Figure Legends**

#### 415 **Figure 1. Inferring Differentiation Trajectories at organismal scale**

416 A. Pijuan-Sala et al. (2019) layout containing single-cell transcriptomes belonging from E6.5 to  
417 E8.5, colored by sampled time-point (left) and by cell-type (right). The overlaying arrows result  
418 from applying the scVelo pipeline to the whole embryonic dataset and represent inferred  
419 developmental trajectories. Arrowheads highlight the erythroid branch, displaying scVelo  
420 trajectory predictions that are inconsistent with real-time sampling.

421 B. Pijuan-Sala et al. (2019) layout highlighting single-cell transcriptomes belonging to E7.5 (left)  
422 and E8.5 (right) and colored by cell-type (see legend in A). The overlaying arrows result from  
423 applying the scVelo pipeline to these individual time-points and represent inferred  
424 developmental trajectories.

#### 425 **Figure 2. Unspliced counts contribute to explaining the variability among cell types**

426 A. Dimensionality reduction with the first two principal components/MOFA factors using spliced  
427 reads alone (left), unspliced reads alone (middle) and both spliced and unspliced (right).  
428 Single-cell transcriptomes are colored by cell-type annotation; see Figure 1 for full legend.

429 B. MOFA characterization of spliced and unspliced reads assessing proportion of variance  
430 explained (i), overlap in highly variable genes calculating using either spliced or unspliced  
431 reads (ii), and factor weight distributions (iii).

#### 432 **Figure 3. A set of genes with complex expression kinetics confounds velocity estimation in 433 erythropoiesis**

434 A. Illustration of phase plot representation in datasets of differentiating cell populations, and  
435 associated scVelo predictions

436 B. Illustration of strategy for MURK gene identification

437 C. Phase plots of representative MURK genes. X-axis: normalized imputed counts of spliced  
438 transcript; y-axis: normalized imputed counts of unspliced transcript.

439 D. GO-term enrichment of MURK genes identified in mouse yolk sac erythropoiesis

440 E. Zoomed-in UMAP of the erythroid branch (see Figure 1 for full UMAP) with scVelo  
441 calculations, before and after removing MURK genes identified in B.

442 **Figure 4. *In vivo* analysis of Gata1 function using a chimaera assay coupled with scRNA-Seq**

- 443 A. Schematic of the G1ER system (Tsang et al., 1997; Weiss et al., 1997)  
444 B. Behaviour of the 73 MURK genes identified in Figure 3 upon Gata1 induction in the G1ER  
445 system (Wu et al., 2011). Wu et al. report that upon Gata1 induction they obtained a total of  
446 2769 upregulated genes, 6079 mildly upregulated, 3566 downregulated, and 3445 with no  
447 response.  
448 C. UMAPS of *Gata1*<sup>-</sup> chimera cells allocated a hemato-endothelial identity colored by cell-type  
449 (sub-clusters defined in Pijuan-Sala et al. (2019) - BP: Blood Progenitors, EC: Endothelial Cells,  
450 Haem: Hemato-endothelial Progenitors, Mk: Megakaryocytes, My: Myeloid cells, Ery:  
451 Erythroid cells) and split by genotype.  
452 D. UMAPS of *Gata1*<sup>-</sup> chimera cells allocated a hemato-endothelial identity colored by sampling  
453 timepoint and split by genotype.  
454 E. Barplots with the quantification of chimera cells mapping to each hemato-endothelial lineage  
455 of the reference dataset (left) and to sampled time-points of the reference dataset (right).

456 **Figure 5. Gata1 chimaera assay reveals disruption of MURK genes and perturbed yolk sac**  
457 **hematopoiesis**

- 458 A. Violin plots of representative genes differentially regulated in *Gata1*<sup>-</sup> hematopoietic lineages.  
459 B. GO-term enrichment of genes downregulated in *Gata1*<sup>-</sup> Ery1 cells compared to their WT  
460 counterparts in chimeras.  
461 C. Venn diagram showing overlap between MURK genes and genes downregulated in *Gata1*<sup>-</sup>  
462 Ery1 cells  
463 D. Phase plots of MURK genes identified along erythroid differentiation, in E8.5 *Gata1*<sup>-</sup> chimera  
464 datasets, colored by tdTom status.

466 **Figure 6. Concept of dual kinetics of gene expression is also revealed in human foetal liver**  
467 **hematopoiesis**

- 468 A. UMAP representation of human fetal liver erythroid cell populations. The overlaying arrows  
469 result from applying the scVelo pipeline using all genes (left) or after MURK gene exclusion  
470 (right), and represent inferred developmental trajectories. Bottom UMAPs are colored by  
471 corresponding scVelo-inferred latent time.

- 472 B. Phase plots of representative MURK genes identified in human fetal liver erythropoiesis  
473 single-cell RNAseq dataset. See (A) for colour code.  
474 C. GO-term enrichment of MURK genes identified in human fetal liver erythropoiesis.

#### 475 **Supplementary Figures**

- 476 1. Dimensionality reduction with the first three principal components/MOFA factors using  
477 spliced reads alone (left), unspliced reads alone (middle) and both spliced and unspliced  
478 (right). Single-cell transcriptomes are colored by cell-type annotation; see Figure 1 for full  
479 legend.
- 480 2. Identification of MURK genes along yolk sac erythropoiesis. A. Phase plots of representative  
481 scVelo driver genes, with scVelo model prediction overlaid (see also Supplementary Table  
482 1). B. Gene ontology of the 73 MURK genes identified in mouse yolk sac erythropoiesis,  
483 including sub-categories. C. Distribution of annotated cell type (top) and sampling time-point  
484 (bottom) along scVelo calculated latent time, using all genes (left panels) and after removing  
485 the MURK genes identified in Figure 3B-C.
- 486 3. Pijuan-Sala et al. (2019) layout highlighting nearest neighbours of *Gata1*<sup>-</sup> chimeras. In red are  
487 nearest neighbours of tdTom<sup>+</sup> mutant cells, in black those of tdTom<sup>-</sup> wildtype cells. To  
488 compare with Figure 1A.
- 489 4. Impact of *Gata1* knockout on *Spi1*/PU.1 expression on the hematoendothelial cell types. X-  
490 axis: *Spi1* log<sub>2</sub>(fold-change) in *Gata1*<sup>-</sup> vs WT chimera cells and Atlas nearest neighbours. Y-axis:  
491 log<sub>10</sub>(FDR).

492

#### 493 **Supplementary Tables**

- 494 1. Driver genes of the scVelo predictions along erythroid differentiation, ranked by likelihood in  
495 the dynamic model.
- 496 2. List of MURK genes identified in Figure 3B-C, ranked by calculated increase in slope value.
- 497 3. Differential Expression Analysis of *Gata1*<sup>-</sup> tdTom<sup>+</sup> vs WT tdTom<sup>-</sup> chimera cells. For the Mk  
498 subset, given the low numbers of WT chimera cells present, the nearest neighbors from the  
499 reference Atlas dataset were included in the comparison. LFC: log fold change.

500

501

502

## 503 **Methods**

### 504 **scVelo implementation**

505 **Mouse atlas dataset.** To obtain separated count matrices for spliced and unspliced mRNAs, we ran  
506 velocity 0.17.17 (La Manno et al., 2018) on the .bam files from the mouse atlas in Pijuan-Sala et al.  
507 2019 (Pijuan-Sala et al., 2019; GEO accession number: GSE87038). We kept all cells that passed the  
508 QC as described in the original publication, but filtered out from downstream analysis the  
509 extraembryonic tissues: ExE endoderm, ExE ectoderm and Parietal endoderm as well as samples with  
510 no timepoint allocation (labelled as 'mixed gastrulation'). Then, we applied scVelo's pipeline v0.2.1  
511 (Bergen et al., 2020). That is, we removed genes with less than 20 shared counts between spliced and  
512 unspliced counts, before normalising and log transforming the remaining genes. Then, we selected  
513 the top 5,000 highly variable genes (HVGs) for further calculation of moments; while performing  
514 imputation using the top 30 nearest neighbours from the graph connectivities generated with the  
515 original UMAP coordinates from Pijuan-Sala et al. 2019. The velocity vectors were computed in  
516 dynamical mode rather than steady state.

517 **Human dataset.** We first downloaded raw reads from Popescu et al., 2019 (Popescu et al., 2019; GEO  
518 accession number: GSE127980), and aligned them against the human genome hg19-3.0.0 with  
519 Cell Ranger v3.0.2 to generate the .bam files and obtain separated count matrices for spliced and  
520 unspliced mRNAs as described above. We filtered out cells with less than 3,550 counts, less than 900  
521 genes and more than 6% mitochondrial counts. We ran scVelo's pipeline selecting 1,500 HVGs to  
522 compute PCA coordinates and applied batch correction using the function reducedMNN from the  
523 batchelor package v1.4.0 (Haghverdi et al., 2018), followed by the estimation of velocity vectors in the  
524 same way it was done for the mouse atlas.

525

### 526 **MOFA+ implementation**

527 We ran MOFA+ v1.4.0 (Argelaguet et al., 2020) using as input the two single cell experiment objects  
528 obtained from the spliced and unspliced counts independently. Each object was created in R using the  
529 scran v1.16.0 (Lun et al., 2016) library as follows: we started from the raw counts, normalized them  
530 with factor sizes obtained after pre-clustering, log transformed and reduced to 5000 HVG. We then  
531 switched to Python v3.7.4, where we regressed out the sample effect and scaled the object to  
532 generate a MOFA+ model with standard parameters. Finally, we used reducedMNN to correct the  
533 MOFA Factors for batch effects. The same objects used as MOFA input were used for PCA calculation  
534 in Figure 2A.

535

### 536 **MURK genes identification**

537 To identify MURK genes, we considered the imputed counts resulting from the scVelo standard  
538 pipeline. Then, for each gene and each population among the Erythroid lineage, we calculated the  
539 unspliced versus spliced slope with a linear regression, as well as the standard error on the slope. We  
540 selected all genes for which the slope in Erythroid3 is significantly higher than the slope in Erythroid2  
541 (according to a one-sided t-test p-value < 0.05), the average spliced counts in Erythroid3 is higher than  
542 the average spliced counts in every other population, and the slope in Erythroid3 positive. We found  
543 73 genes that respect all these criteria in the mouse dataset, 377 in the human dataset.

544

### 545 **Gene ontology enrichment analysis**

546 We performed gene ontology enrichment analysis using the <http://geneontology.org> website  
547 comparing the MURK genes against all biological processes, with the default all *Mus musculus* genes  
548 in database as background set (Ashburner et al., 2000; The Gene Ontology, 2019). We ranked the  
549 processes by FDR.

550

### 551 **Overlap testing**

552 Overlap was tested with Fisher exact test. We calculated the probability of having  $m = 45$  genes of our  
553  $n = 73$  MURK genes mapping to the  $A = 1022$  high response genes (out of  $N = 4195$  genes) in the Wu  
554 et al., 2011 publication (GEO accession number: GSE30142) as the probability of randomly picking  $m$   
555 elements of a specific type when randomly choosing  $n$  elements out of  $N$ , where the frequency of the  
556 special type is  $A/N$ .

557

### 558 **Gata1<sup>-</sup> chimera dataset generation and analysis**

559 **Embryo collection.** All procedures were performed in strict accordance to the UK Home Office  
560 regulations for animal research under the project license number PPL 70/8406. **Chimaera generation.**  
561 TdTomato-expressing mouse embryonic stem cells (ESC) were derived as previously described (Pijuan-  
562 Sala et al., 2019). Briefly, ESC lines were derived from E3.5 blastocysts obtained by crossing a male  
563 ROSA26tdTomato (Jax Labs – 007905) with a wildtype C57BL/6 female, expanded under the 2i+LIF  
564 conditions (Ying et al., 2008) and transiently transfected with a Cre-IRES-GFP plasmid (Wray et al.,

565 2011) using Lipofectamine 3000 Transfection Reagent (ThermoFisher Scientific, #L3000008) according  
566 to manufacturer's instructions. A tdTomato-positive, male, karyotypically normal line, competent for  
567 chimaera generation as assessed using morula aggregation assay, was selected for targeting *Gata1*.  
568 Two guides were designed using the <http://crispr.mit.edu> tool (guide 1: CGGCTACTCCACTGTGGCGG;  
569 guide 2: CGCTTCTGGGCCGGATGAG) and were cloned into the pX458 plasmid (Addgene, #48138) as  
570 previously described (Ran et al., 2013). The obtained plasmids were then used to transfect the cells  
571 and single transfected clones were expanded and assessed for Cas9-induced mutations. Genomic DNA  
572 was isolated by incubating cell pellets in 0.1 mg/ml of Proteinase K (Sigma, #03115828001) in TE buffer  
573 at 50°C for 2 hours, followed by 5 min at 99°C. The sequence flanking the guide-targeted sites was  
574 amplified from the genomic DNA by polymerase chain reaction (PCR) in a Biometra T3000  
575 Thermocycler (30 sec at 98°C ; 30 cycles of 10 sec at 98°C, 20 sec at 58°C, 20 sec at 72°C; and  
576 elongation for 7 min at 72°C) using the Phusion High-Fidelity DNA Polymerase (NEB, #M0530S)  
577 according to the manufacturer's instructions. Primers including Nextera overhangs were used (F-  
578 TCGTCGGCAGCGTCAGATGTGTATAAGAGACAGTCTACCTGCCTCAACTGTG; R-  
579 GTCTCGTGGCTCGGAGATGTGTATAAGAGACAGTCTGTCTTGGGCAGGAACA), allowing library  
580 preparation with the Nextera XT Kit (Illumina, #15052163), and sequencing was performed using the  
581 Illumina MiSeq system according to manufacturer's instructions. An ESC clone showing a 38 base-pair  
582 frameshift mutation in exon 4 resulting in the functional inactivation of *Gata1* were selected for  
583 injection into C57BL/6 E3.5 blastocysts. A total of 6 chimaeric embryos were harvested at E8.5,  
584 dissected, and single-cell suspensions were generated from three independent pools of embryos by  
585 TrypLE Express dissociation reagent (Thermo Fisher Scientific) incubation for 7-10 minutes at 37°C  
586 under agitation. Single-cell suspensions were sorted into tdTom+ and tdTom- samples using a BD Influx  
587 sorter with DAPI at 1µg/ml (Sigma) as a viability stain for subsequent 10X scRNA-seq library  
588 preparation (version 3 chemistry), and sequencing using an S1 flow cell in the Illumina Novaseq  
589 platform, which resulted in 8420 tdTom<sup>-</sup> and 7944 tdTom<sup>+</sup> cells that passed quality control (see  
590 "Single-cell RNA sequencing analysis" below).

591 **Single-cell RNA sequencing analysis.** Raw files were processed with Cell Ranger 3.0.2 using default  
592 mapping arguments. Reads were mapped to the mm10 genome and counted with GRCm38.92  
593 annotation, including tdTomato sequence for chimera cells. Cell barcodes with expression profiles  
594 significantly different to the ambient mRNA expression profile were identified using emptyDrops (Lun  
595 et al., 2019), and cell barcodes with low complexity, i.e. low total mRNA counts and/or high  
596 mitochondrial proportion, were identified by fitting four-component bivariate mixture models to the  
597 log<sub>10</sub>-transformed total mRNA counts and percentage of mitochondrial counts, and selecting the  
598 components with high total mRNA and low mitochondrial percentage. Gene expression normalization

599 and doublet cell barcodes were identified using the approach taken by Pijuan-Sala et al. (2019). Both  
600 spliced and unspliced count matrices were extracted using velocity 0.17.17 (La Manno et al., 2018).

601 **Mapping to the reference dataset.** We mapped the chimaera cells to the mouse atlas following almost  
602 exactly the procedure used in the original publication article to map the *Tal1* chimaera. First, we  
603 concatenated the mouse atlas and chimaera counts (both previously controlled for quality of the  
604 cells), normalized the resulting counts matrix with scran, computed HVGs and then applied  
605 multiBatchPCA, and reducedMNN with cosine normalization from batchelor (Haghverdi et al., 2018)  
606 for batch effect correction within samples (where sample refers to a single lane of a 10x Chromium  
607 chip) as well as between datasets in order to extract a number of nearest neighbours between the  
608 mouse atlas and the chimaera using queryKNN from BiocNeighbors package v1.6.0.

609 **Differential Gene Expression Analysis.** For differential gene expression analysis, we took samples that  
610 included at least 7 cells per tdTom status per cell population (eg. Erythroid3). We ran the analysis in  
611 scanpy v1.5.1 (Wolf et al., 2018) with Wilcoxon test and choosing 2 as fold change and 0.1 as false  
612 discovery rate thresholds.

613

614

## 615 **References**

616 Akunuru, S., and Geiger, H. (2016). Aging, Clonality, and Rejuvenation of Hematopoietic Stem Cells.  
617 *Trends Mol Med* 22, 701-712.

618 Argelaguet, R., Arnol, D., Bredikhin, D., Deloro, Y., Velten, B., Marioni, J.C., and Stegle, O. (2020).  
619 MOFA+: a statistical framework for comprehensive integration of multi-modal single-cell data.  
620 *Genome Biol* 21, 111.

621 Ashburner, M., Ball, C.A., Blake, J.A., Botstein, D., Butler, H., Cherry, J.M., Davis, A.P., Dolinski, K.,  
622 Dwight, S.S., Eppig, J.T., *et al.* (2000). Gene ontology: tool for the unification of biology. The Gene  
623 Ontology Consortium. *Nat Genet* 25, 25-29.

624 Bergen, V., Lange, M., Peidli, S., Wolf, F.A., and Theis, F.J. (2020). Generalizing RNA velocity to  
625 transient cell states through dynamical modeling. *Nat Biotechnol*.

626 Bhatnagar, N., Nizery, L., Tunstall, O., Vyas, P., and Roberts, I. (2016). Transient Abnormal  
627 Myelopoiesis and AML in Down Syndrome: an Update. *Curr Hematol Malig Rep* 11, 333-341.

628 Borrett, M.J., Innes, B.T., Jeong, D., Tahmasian, N., Storer, M.A., Bader, G.D., Kaplan, D.R., and Miller,  
629 F.D. (2020). Single-Cell Profiling Shows Murine Forebrain Neural Stem Cells Reacquire a  
630 Developmental State when Activated for Adult Neurogenesis. *Cell Rep* 32, 108022.

631 Carotta, S., Wu, L., and Nutt, S.L. (2010). Surprising new roles for PU.1 in the adaptive immune  
632 response. *Immunol Rev* 238, 63-75.

- 633 Choe, K.S., Ujhelly, O., Wontakal, S.N., and Skoultchi, A.I. (2010). PU.1 directly regulates cdk6 gene  
634 expression, linking the cell proliferation and differentiation programs in erythroid cells. *J Biol Chem*  
635 *285*, 3044-3052.
- 636 Dahlin, J.S., Hamey, F.K., Pijuan-Sala, B., Shepherd, M., Lau, W.W.Y., Nestorowa, S., Weinreb, C.,  
637 Wolock, S., Hannah, R., Diamanti, E., *et al.* (2018). A single-cell hematopoietic landscape resolves 8  
638 lineage trajectories and defects in Kit mutant mice. *Blood* *131*, e1-e11.
- 639 Ezer, D., Moignard, V., Gottgens, B., and Adryan, B. (2016). Determining Physical Mechanisms of  
640 Gene Expression Regulation from Single Cell Gene Expression Data. *PLoS Comput Biol* *12*, e1005072.
- 641 Fujiwara, Y., Browne, C.P., Cunniff, K., Goff, S.C., and Orkin, S.H. (1996). Arrested development of  
642 embryonic red cell precursors in mouse embryos lacking transcription factor GATA-1. *Proc Natl Acad*  
643 *Sci U S A* *93*, 12355-12358.
- 644 Gerber, T., Murawala, P., Knapp, D., Masselink, W., Schuez, M., Hermann, S., Gac-Santel, M.,  
645 Nowoshilow, S., Kageyama, J., Khattak, S., *et al.* (2018). Single-cell analysis uncovers convergence of  
646 cell identities during axolotl limb regeneration. *Science* *362*.
- 647 Grosveld, F., van Assendelft, G.B., Greaves, D.R., and Kollias, G. (1987). Position-independent, high-  
648 level expression of the human beta-globin gene in transgenic mice. *Cell* *51*, 975-985.
- 649 Gutierrez, L., Tsukamoto, S., Suzuki, M., Yamamoto-Mukai, H., Yamamoto, M., Philipsen, S., and  
650 Ohneda, K. (2008). Ablation of Gata1 in adult mice results in aplastic crisis, revealing its essential role  
651 in steady-state and stress erythropoiesis. *Blood* *111*, 4375-4385.
- 652 Haghverdi, L., Lun, A.T.L., Morgan, M.D., and Marioni, J.C. (2018). Batch effects in single-cell RNA-  
653 sequencing data are corrected by matching mutual nearest neighbors. *Nat Biotechnol* *36*, 421-427.
- 654 Higgs, D.R., Wood, W.G., Jarman, A.P., Sharpe, J., Lida, J., Pretorius, I.M., and Ayyub, H. (1990). A  
655 major positive regulatory region located far upstream of the human alpha-globin gene locus. *Genes*  
656 *Dev* *4*, 1588-1601.
- 657 Hoppe, P.S., Schwarzfischer, M., Loeffler, D., Kokkaliaris, K.D., Hilsenbeck, O., Moritz, N., Endeke, M.,  
658 Filipczyk, A., Gambardella, A., Ahmed, N., *et al.* (2016). Early myeloid lineage choice is not initiated  
659 by random PU.1 to GATA1 protein ratios. *Nature* *535*, 299-302.
- 660 Kanton, S., Boyle, M.J., He, Z., Santel, M., Weigert, A., Sanchis-Calleja, F., Guijarro, P., Sidow, L., Fleck,  
661 J.S., Han, D., *et al.* (2019). Organoid single-cell genomic atlas uncovers human-specific features of  
662 brain development. *Nature* *574*, 418-422.
- 663 Kim, J.K., and Marioni, J.C. (2013). Inferring the kinetics of stochastic gene expression from single-cell  
664 RNA-sequencing data. *Genome Biol* *14*, R7.
- 665 Kondo, A., Fujiwara, T., Okitsu, Y., Fukuhara, N., Onishi, Y., Nakamura, Y., Sawada, K., and Harigae, H.  
666 (2016). Identification of a novel putative mitochondrial protein FAM210B associated with erythroid  
667 differentiation. *Int J Hematol* *103*, 387-395.
- 668 La Manno, G., Soldatov, R., Zeisel, A., Braun, E., Hochgerner, H., Petukhov, V., Lidschreiber, K.,  
669 Kastrioti, M.E., Lonnerberg, P., Furlan, A., *et al.* (2018). RNA velocity of single cells. *Nature* *560*, 494-  
670 498.

- 671 Li, Z., Godinho, F.J., Klusmann, J.H., Garriga-Canut, M., Yu, C., and Orkin, S.H. (2005). Developmental  
672 stage-selective effect of somatically mutated leukemogenic transcription factor GATA1. *Nat Genet*  
673 *37*, 613-619.
- 674 Liu, X., Chen, Y., Zhang, Y., Liu, Y., Liu, N., Botten, G.A., Cao, H., Orkin, S.H., Zhang, M.Q., and Xu, J.  
675 (2020). Multiplexed capture of spatial configuration and temporal dynamics of locus-specific 3D  
676 chromatin by biotinylated dCas9. *Genome Biol* *21*, 59.
- 677 Lun, A.T., McCarthy, D.J., and Marioni, J.C. (2016). A step-by-step workflow for low-level analysis of  
678 single-cell RNA-seq data with Bioconductor. *F1000Res* *5*, 2122.
- 679 Lun, A.T.L., Riesenfeld, S., Andrews, T., Dao, T.P., Gomes, T., participants in the 1st Human Cell Atlas,  
680 J., and Marioni, J.C. (2019). EmptyDrops: distinguishing cells from empty droplets in droplet-based  
681 single-cell RNA sequencing data. *Genome Biol* *20*, 63.
- 682 Mahadevaiah, S.K., Sangrithi, M.N., Hirota, T., and Turner, J.M.A. (2020). A single-cell transcriptome  
683 atlas of marsupial embryogenesis and X inactivation. *Nature* *586*, 612-617.
- 684 McGrath, K., and Palis, J. (2008). Ontogeny of erythropoiesis in the mammalian embryo. *Curr Top*  
685 *Dev Biol* *82*, 1-22.
- 686 Mettananda, S., Gibbons, R.J., and Higgs, D.R. (2016). Understanding alpha-globin gene regulation  
687 and implications for the treatment of beta-thalassemia. *Ann N Y Acad Sci* *1368*, 16-24.
- 688 Monteiro, R., Pouget, C., and Patient, R. (2011). The *gata1/pu.1* lineage fate paradigm varies  
689 between blood populations and is modulated by *tif1gamma*. *EMBO J* *30*, 1093-1103.
- 690 Palis, J. (2016). Hematopoietic stem cell-independent hematopoiesis: emergence of erythroid,  
691 megakaryocyte, and myeloid potential in the mammalian embryo. *FEBS Lett* *590*, 3965-3974.
- 692 Pevny, L., Lin, C.S., D'Agati, V., Simon, M.C., Orkin, S.H., and Costantini, F. (1995). Development of  
693 hematopoietic cells lacking transcription factor GATA-1. *Development* *121*, 163-172.
- 694 Pevny, L., Simon, M.C., Robertson, E., Klein, W.H., Tsai, S.F., D'Agati, V., Orkin, S.H., and Costantini, F.  
695 (1991). Erythroid differentiation in chimaeric mice blocked by a targeted mutation in the gene for  
696 transcription factor GATA-1. *Nature* *349*, 257-260.
- 697 Pijuan-Sala, B., Griffiths, J.A., Guibentif, C., Hiscock, T.W., Jawaid, W., Calero-Nieto, F.J., Mulas, C.,  
698 Ibarra-Soria, X., Tyser, R.C.V., Ho, D.L.L., *et al.* (2019). A single-cell molecular map of mouse  
699 gastrulation and early organogenesis. *Nature* *566*, 490-495.
- 700 Pina, C., May, G., Soneji, S., Hong, D., and Enver, T. (2008). MLLT3 regulates early human erythroid  
701 and megakaryocytic cell fate. *Cell Stem Cell* *2*, 264-273.
- 702 Popescu, D.M., Botting, R.A., Stephenson, E., Green, K., Webb, S., Jardine, L., Calderbank, E.F.,  
703 Polanski, K., Goh, I., Efremova, M., *et al.* (2019). Decoding human fetal liver haematopoiesis. *Nature*  
704 *574*, 365-371.
- 705 Ran, F.A., Hsu, P.D., Wright, J., Agarwala, V., Scott, D.A., and Zhang, F. (2013). Genome engineering  
706 using the CRISPR-Cas9 system. *Nat Protoc* *8*, 2281-2308.

- 707 Roberts, I., Alford, K., Hall, G., Juban, G., Richmond, H., Norton, A., Vallance, G., Perkins, K., Marchi,  
708 E., McGowan, S., *et al.* (2013). GATA1-mutant clones are frequent and often unsuspected in babies  
709 with Down syndrome: identification of a population at risk of leukemia. *Blood* *122*, 3908-3917.
- 710 Schultz, M.B., and Sinclair, D.A. (2016). When stem cells grow old: phenotypes and mechanisms of  
711 stem cell aging. *Development* *143*, 3-14.
- 712 Shivdasani, R.A., Fujiwara, Y., McDevitt, M.A., and Orkin, S.H. (1997). A lineage-selective knockout  
713 establishes the critical role of transcription factor GATA-1 in megakaryocyte growth and platelet  
714 development. *EMBO J* *16*, 3965-3973.
- 715 Storry, J.R., Joud, M., Christophersen, M.K., Thuresson, B., Akerstrom, B., Sojka, B.N., Nilsson, B., and  
716 Olsson, M.L. (2013). Homozygosity for a null allele of SMIM1 defines the Vel-negative blood group  
717 phenotype. *Nat Genet* *45*, 537-541.
- 718 The Gene Ontology, C. (2019). The Gene Ontology Resource: 20 years and still GOing strong. *Nucleic  
719 Acids Res* *47*, D330-D338.
- 720 Ton, M.N., Guibentif, C., and Gottgens, B. (2020). Single cell genomics and developmental biology:  
721 moving beyond the generation of cell type catalogues. *Curr Opin Genet Dev* *64*, 66-71.
- 722 Tsang, A.P., Visvader, J.E., Turner, C.A., Fujiwara, Y., Yu, C., Weiss, M.J., Crossley, M., and Orkin, S.H.  
723 (1997). FOG, a multitype zinc finger protein, acts as a cofactor for transcription factor GATA-1 in  
724 erythroid and megakaryocytic differentiation. *Cell* *90*, 109-119.
- 725 Vyas, P., Ault, K., Jackson, C.W., Orkin, S.H., and Shivdasani, R.A. (1999). Consequences of GATA-1  
726 deficiency in megakaryocytes and platelets. *Blood* *93*, 2867-2875.
- 727 Wagner, D.E., Weinreb, C., Collins, Z.M., Briggs, J.A., Megason, S.G., and Klein, A.M. (2018). Single-  
728 cell mapping of gene expression landscapes and lineage in the zebrafish embryo. *Science* *360*, 981-  
729 987.
- 730 Weinreb, C., Rodriguez-Fraticelli, A., Camargo, F.D., and Klein, A.M. (2020). Lineage tracing on  
731 transcriptional landscapes links state to fate during differentiation. *Science* *367*.
- 732 Weiss, M.J., Yu, C., and Orkin, S.H. (1997). Erythroid-cell-specific properties of transcription factor  
733 GATA-1 revealed by phenotypic rescue of a gene-targeted cell line. *Mol Cell Biol* *17*, 1642-1651.
- 734 Wolf, F.A., Angerer, P., and Theis, F.J. (2018). SCANPY: large-scale single-cell gene expression data  
735 analysis. *Genome Biol* *19*, 15.
- 736 Wray, J., Kalkan, T., Gomez-Lopez, S., Eckardt, D., Cook, A., Kemler, R., and Smith, A. (2011).  
737 Inhibition of glycogen synthase kinase-3 alleviates Tcf3 repression of the pluripotency network and  
738 increases embryonic stem cell resistance to differentiation. *Nat Cell Biol* *13*, 838-845.
- 739 Wu, W., Cheng, Y., Keller, C.A., Ernst, J., Kumar, S.A., Mishra, T., Morrissey, C., Dorman, C.M., Chen,  
740 K.B., Drautz, D., *et al.* (2011). Dynamics of the epigenetic landscape during erythroid differentiation  
741 after GATA1 restoration. *Genome Res* *21*, 1659-1671.
- 742 Ying, Q.L., Wray, J., Nichols, J., Batlle-Morera, L., Doble, B., Woodgett, J., Cohen, P., and Smith, A.  
743 (2008). The ground state of embryonic stem cell self-renewal. *Nature* *453*, 519-523.

- 744 Zhang, P., Behre, G., Pan, J., Iwama, A., Wara-Aswapati, N., Radomska, H.S., Auron, P.E., Tenen, D.G.,  
745 and Sun, Z. (1999). Negative cross-talk between hematopoietic regulators: GATA proteins repress  
746 PU.1. *Proc Natl Acad Sci U S A* *96*, 8705-8710.
- 747 Zhang, Q., He, Y., Luo, N., Patel, S.J., Han, Y., Gao, R., Modak, M., Carotta, S., Haslinger, C., Kind, D.,  
748 *et al.* (2019). Landscape and Dynamics of Single Immune Cells in Hepatocellular Carcinoma. *Cell* *179*,  
749 829-845 e820.
- 750 Zhou, W., Yui, M.A., Williams, B.A., Yun, J., Wold, B.J., Cai, L., and Rothenberg, E.V. (2019). Single-Cell  
751 Analysis Reveals Regulatory Gene Expression Dynamics Leading to Lineage Commitment in Early T  
752 Cell Development. *Cell Syst* *9*, 321-337 e329.

# Figure 1: Inferring Differentiation Trajectories at organismal scale

bioRxiv preprint doi: <https://doi.org/10.1101/2020.12.21.423773>; this version posted December 22, 2020. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under a [CC-BY-NC-ND 4.0 International license](https://creativecommons.org/licenses/by-nc-nd/4.0/).

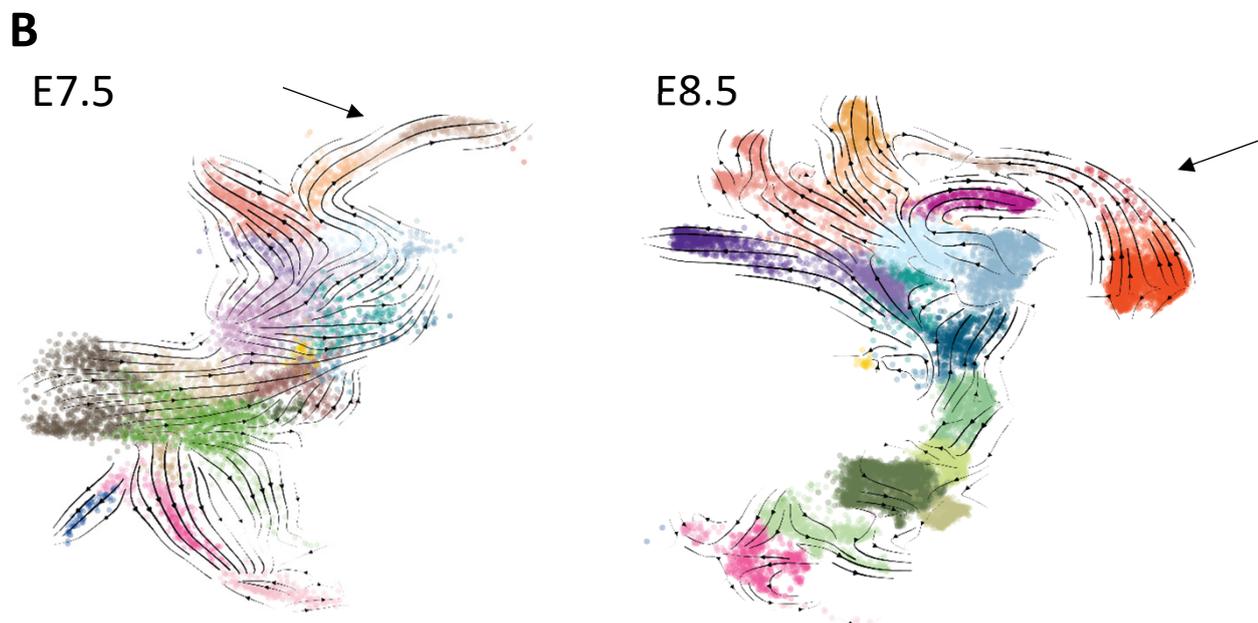
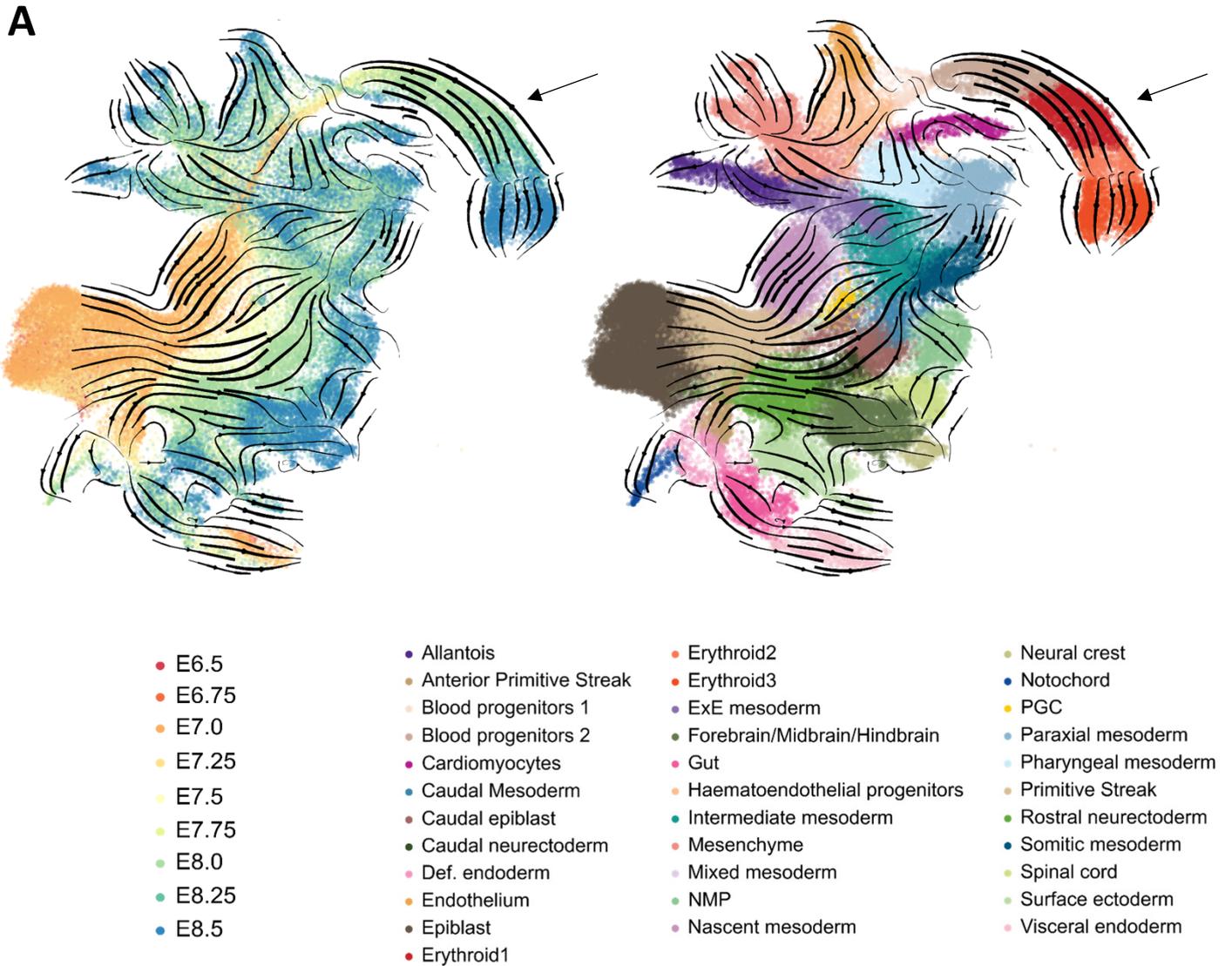


Figure 2: Unspliced counts contribute to explaining the variability among cell types

bioRxiv preprint doi: <https://doi.org/10.1101/2020.12.21.423773>; this version posted December 22, 2020. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under a [CC-BY-NC-ND 4.0 International license](#).

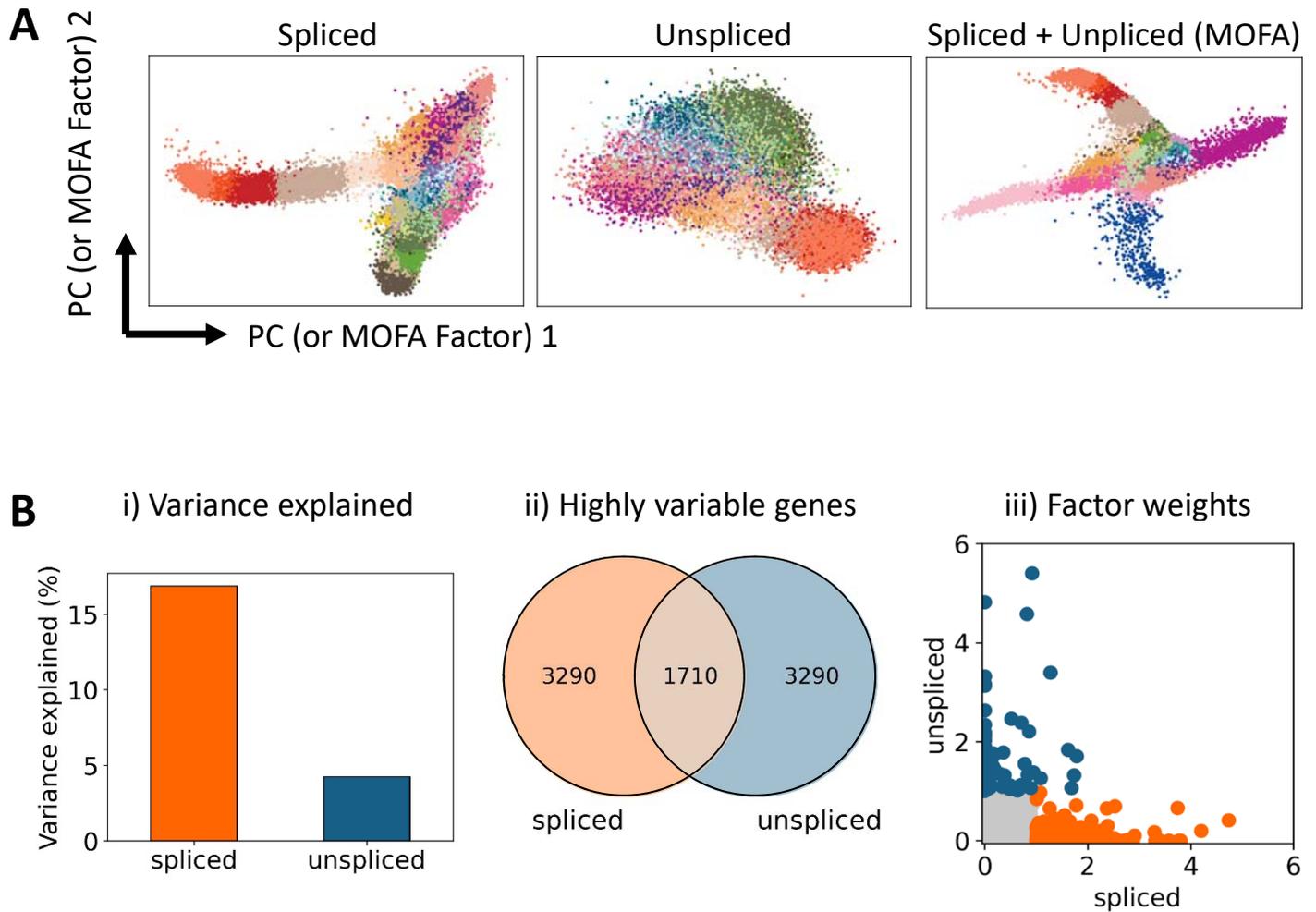


Figure 3. A set of genes with complex expression kinetics confounds velocity estimation in erythropoiesis

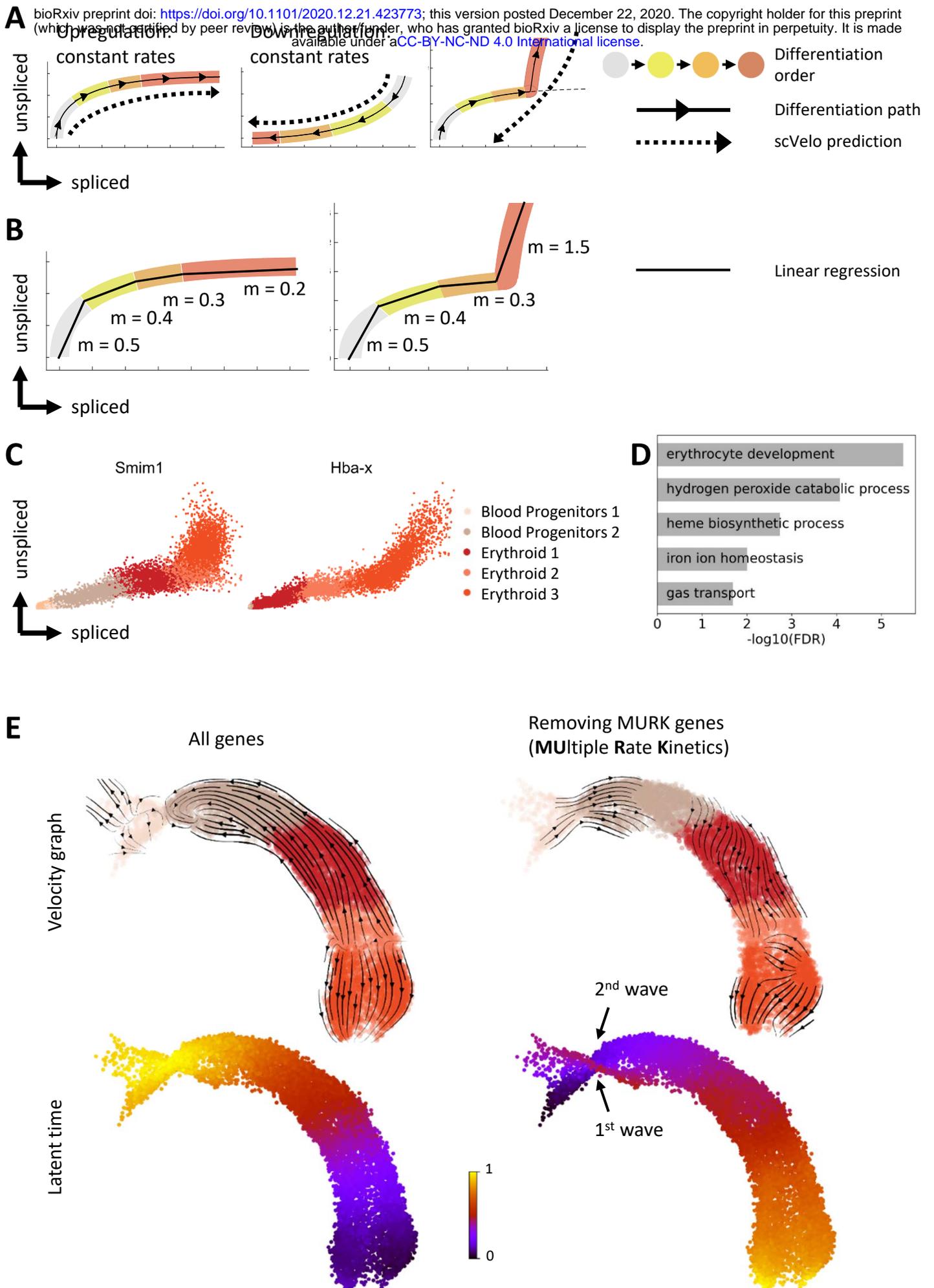
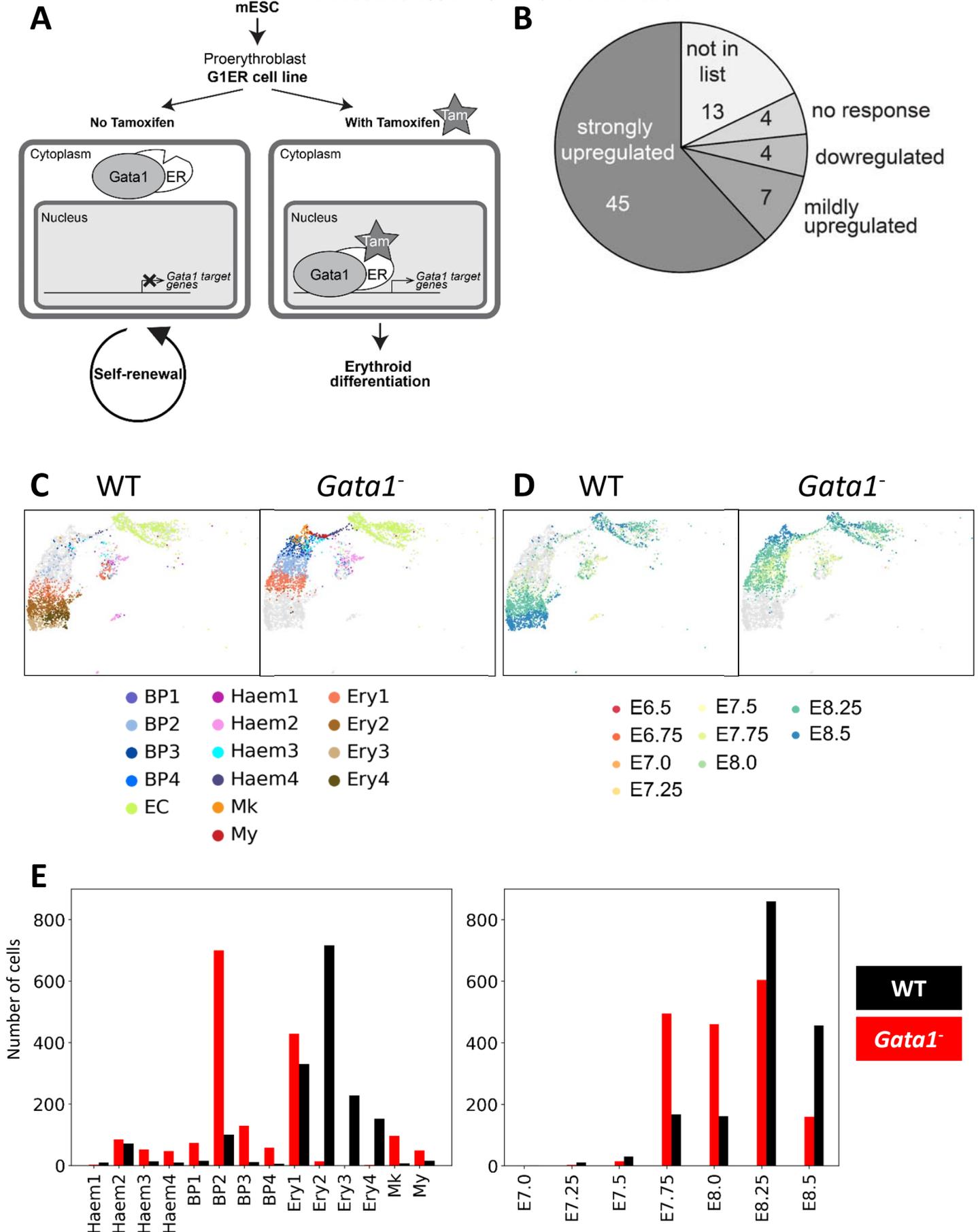


Figure 4: In vivo analysis of Gata1 function using Chimaera assays coupled with scRNA-Seq

bioRxiv preprint doi: <https://doi.org/10.1101/2020.12.21.423773>; this version posted December 22, 2020. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under a [CC-BY-NC-ND 4.0 International license](#).



## Figure 5. Gata1 chimaera assay reveals disruption of MURK genes and perturbed yolk sac hematopoiesis

bioRxiv preprint doi: <https://doi.org/10.1101/2020.12.21.423773>; this version posted December 22, 2020. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under a [CC-BY-NC-ND 4.0 International license](https://creativecommons.org/licenses/by-nc-nd/4.0/).

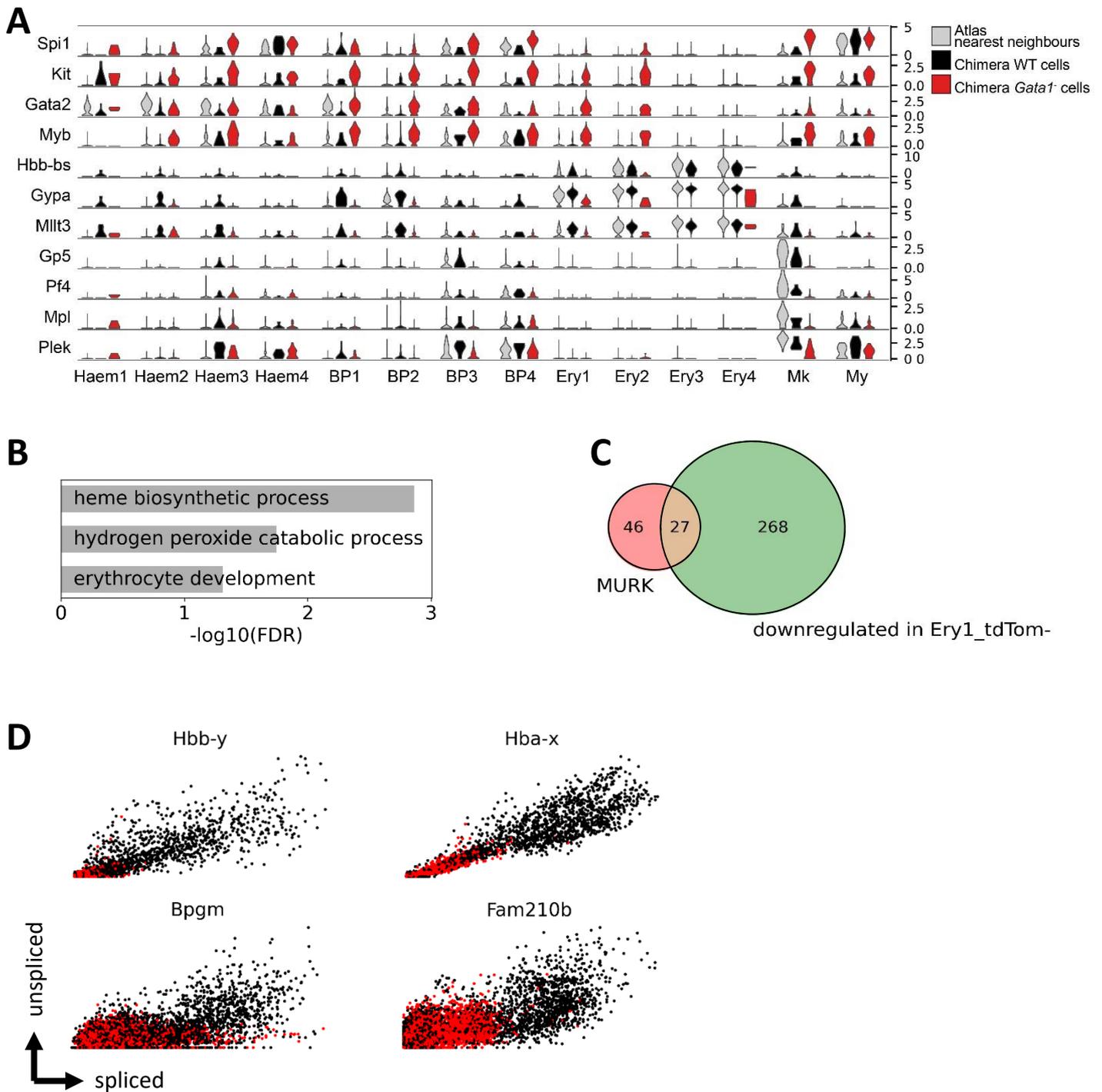


Figure 6. Concept of dual kinetics of gene expression is also revealed in human foetal liver hematopoiesis

**A** bioRxiv preprint doi: <https://doi.org/10.1101/2020.12.21.423773>; this version posted December 22, 2020. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under a [CC-BY-NC-ND 4.0 International license](https://creativecommons.org/licenses/by-nc-nd/4.0/).

