



THE UNIVERSITY *of* EDINBURGH

Edinburgh Research Explorer

Single cell analyses and machine learning define hematopoietic progenitor and HSC-like cells derived from human PSCs

Citation for published version:

Fidanza, A, Stumpf, PS, Ramachandran, P, Tamagno, S, Babbie, A, Lopez Yrigoyen, M, Taylor, H, Easterbrook, J, Henderson, B, Axton, R, Henderson, NC, Medvinsky, A, Ottersbach, K, Romanò, N & Forrester, L 2020, 'Single cell analyses and machine learning define hematopoietic progenitor and HSC-like cells derived from human PSCs', *Blood*, vol. 136, no. 25, pp. 2893–2904.
<https://doi.org/10.1182/blood.2020006229>

Digital Object Identifier (DOI):

[10.1182/blood.2020006229](https://doi.org/10.1182/blood.2020006229)

Link:

[Link to publication record in Edinburgh Research Explorer](#)

Document Version:

Peer reviewed version

Published In:

Blood

General rights

Copyright for the publications made accessible via the Edinburgh Research Explorer is retained by the author(s) and / or other copyright owners and it is a condition of accessing these publications that users recognise and abide by the legal requirements associated with these rights.

Take down policy

The University of Edinburgh has made every reasonable effort to ensure that Edinburgh Research Explorer content complies with UK legislation. If you believe that the public display of this file breaches copyright please contact openaccess@ed.ac.uk providing details, and we will remove access to the work immediately and investigate your claim.



Single cell analyses and machine learning define hematopoietic progenitor and HSC-like cells derived from human PSCs.

Tracking no: BLD-2020-006229R1

Antonella Fidanza (University of Edinburgh, United Kingdom) Patrick Stumpf (University of Southampton, United Kingdom) Prakash Ramachandran (University of Edinburgh, United Kingdom) Sara Tamagno (University of Edinburgh, United Kingdom) Ann Babbie (Centre for Integrative Systems Biology and Bioinformatics, Department of Life Sciences, Imperial College London, United Kingdom) Martha Lopez-Yrigoyen (University of Edinburgh, United Kingdom) Alice Taylor (University of Edinburgh,) Jennifer Easterbrook (University of Edinburgh, United Kingdom) Beth Henderson (University of Edinburgh, United Kingdom) Richard Axton (University of Edinburgh, United Kingdom) Neil Henderson (University of Edinburgh, United Kingdom) Alexander Medvinsky (The University of Edinburgh, United Kingdom) Katrin Ottersbach (University of Edinburgh, United Kingdom) Nicola Romano (University of Edinburgh, United Kingdom) Lesley Forrester (University of Edinburgh, United Kingdom)

Abstract:

Haematopoietic stem and progenitor cells (HSPCs) develop through distinct waves at various anatomical sites during embryonic development. The *in vitro* differentiation of human pluripotent stem cells (hPSCs) is able to recapitulate some of these processes but it has proven difficult to generate functional haematopoietic stem cells (HSCs). To define the dynamics and heterogeneity of HSPCs that can be generated *in vitro* from hPSCs, we exploited single cell RNA sequencing (scRNAseq) in combination with single cell protein expression analysis. Bioinformatics analyses and functional validation defined the transcriptomes of naïve progenitors as well as erythroid, megakaryocyte and leukocyte-committed progenitors and we identified CD44, CD326, ICAM2/CD9 and CD18 as novel markers of these progenitors, respectively. Using an artificial neural network (ANN), that we trained on a scRNAseq derived from human fetal liver, we were able to identify a wide range of hPSCs-derived HPSC phenotypes, including a small group classified as HSCs. This transient HSC-like population reduced as differentiation proceeded and was completely missing in the dataset that had been generated using cells selected on the basis of CD43 expression. By comparing the single cell transcriptome of *in vitro*-generated HSC-like cells with those generated within the fetal liver we identified transcription factors and molecular pathways that can be targeted with the aim of improving HSC differentiation *in vitro*.

Conflict of interest: No COI declared

COI notes:

Preprint server: Yes; BioRxiv <https://doi.org/10.1101/602565>

Author contributions and disclosures: AF, designed and performed research, analyzed the data and wrote the manuscript. AF, PS, AB and NR performed bioinformatics analysis. PR, ST, MLY, AHT, JE, BH, RA performed research. LMF designed the experiment, analyzed data and wrote the manuscript. NH, AM, KO, and NR provided intellectual input and final approval of the manuscript.

Non-author contributions and disclosures: No;

Agreement to Share Publication-Related Data and Data Sharing Statement: We have created a webpage where the data can be freely browsed, plots can be generated and exported, and full datasets can be downloaded. The link is provided in the manuscript.

Clinical trial registration information (if any):

1 **Single cell multimodal analyses and machine learning define haematopoietic**
2 **progenitor and HSC-like cells derived in vitro from human pluripotent stem cells.**

3
4 Antonella Fidanza^{1*}, Patrick S Stumpf², Prakash Ramachandran³, Sara Tamagno¹, Ann
5 Babbie⁴, Martha Lopez-Yrigoyen¹, A. Helen Taylor¹, Jennifer Easterbrook¹, Beth
6 Henderson³, Richard Axton¹, Neil C. Henderson³, Alexander Medvinsky¹, Katrin
7 Ottersbach¹, Nicola Romanò⁵, Lesley M. Forrester^{1*}.

8
9 1 - Centre for Regenerative Medicine, University of Edinburgh, Edinburgh, UK

10 2 - Joint Research Center for Computational Biomedicine, Uniklinik RWTH Aachen,
11 Aachen, Germany

12 3 - Centre for Inflammation Research, University of Edinburgh, Edinburgh, UK

13 4 - Centre for Integrative Systems Biology and Bioinformatics, Department of Life
14 Sciences, Imperial College London, London, UK

15 5 - Centre for Discovery Brain Sciences, University of Edinburgh, Edinburgh, UK

16
17 *correspondence to afidanza@ed.ac.uk and L.Forrester@ed.ac.uk

18
19 **Key point 1** - Single-cell and CITE-seq profiling of human HSPCs derived in vitro from
20 pluripotent stem cells browsable at <http://188.166.158.65/scRNAseqHPC/>

21 **Key point 2** – Artificial Neural Network identifies HSC-like cells derived in vitro from
22 hPSCs.

23
24 **Abstract**

25 Haematopoietic stem and progenitor cells (HSPCs) develop through distinct waves at
26 various anatomical sites during embryonic development. The in vitro differentiation of
27 human pluripotent stem cells (hPSCs) is able to recapitulate some of these processes,
28 however, it has proven difficult to generate functional haematopoietic stem cells (HSCs).
29 To define the dynamics and heterogeneity of HSPCs that can be generated in vitro from
30 hPSCs, we exploited single cell RNA sequencing (scRNAseq) in combination with single
31 cell protein expression analysis. **Bioinformatics analyses and functional validation defined**
32 **the transcriptomes of naïve progenitors as well as erythroid, megakaryocyte and**
33 **leukocyte-committed progenitors and we identified CD44, CD326, ICAM2/CD9 and CD18**

34 as markers of these progenitors, respectively. Using an artificial neural network (ANN),
35 that we trained on a scRNAseq derived from human fetal liver, we were able to identify a
36 wide range of hPSCs-derived HPSC phenotypes, including a small group classified as
37 HSCs. This transient HSC-like population decreased as differentiation proceeded and was
38 completely missing in the dataset that had been generated using cells selected on the
39 basis of CD43 expression. By comparing the single cell transcriptome of in vitro-generated
40 HSC-like cells with those generated within the fetal liver we identified transcription factors
41 and molecular pathways that can be exploited in the future to improve the in vitro
42 production of HSCs.

43

44 Introduction

45 Human pluripotent stem cells (hPSCs) can be differentiated in vitro into various
46 haematopoietic cell types, providing both a model for basic research studies and a source
47 of clinically relevant cells¹. During embryonic development, two waves of restricted
48 haematopoietic progenitors arise in the extraembryonic tissues of the yolk sac, before
49 emergence of haematopoietic stem cells (HSCs) in the embryo proper². In the mouse
50 embryo, at E7.25 the first “primitive” wave gives rise to erythrocytes, megakaryocytes and
51 macrophages^{3,4}, after at E8.25, the second wave, also known as the first “definitive”
52 progenitors, consists of erythro-myeloid progenitors (EMPs) that can be distinguished from
53 the primitive progenitors by their potential to generate granulocytes⁵. Furthermore, the
54 monocytes that emerge from EMPs provide the embryo with tissue resident macrophage;
55 the first life-long lasting population of immune cells⁶⁻⁸. Intraembryonic hematopoiesis is
56 established during E10.5-E11.5 in the aorta-gonad-mesonephros (AGM) region with the
57 emergence of HSCs, that will sustain the lifespan production of all blood lineages, also
58 upon transplantation⁹. A number of studies have indicated that human haematopoietic
59 development follows a comparable process¹⁰⁻¹³ but for ethical reasons it has proven
60 difficult to gain a clear insight into the lineage potential and hierarchical relationships
61 between early human haematopoietic progenitors. The dynamic nature and the
62 heterogeneity of haematopoietic progenitor populations that arise during development
63 poses additional confounders to the identification of molecular mechanisms associated
64 with their formation and function.

65 To gain insight into the transcriptome of developing human haematopoietic progenitors,
66 we performed in-depth characterization of haematopoietic progenitors derived from

67 hPSCs. Single-cell expression profiles of hPSCs-derived haematopoietic cells have been
68 reported but these previous studies either used a limited number of cells or used biased
69 approaches in their isolation and so failed to depict their trajectory of differentiation^{14–16}.
70 This significantly impacted on the ability to resolve the complex heterogeneity of the
71 progenitor pool, to identify the hierarchal relationship between subpopulations and to
72 compare the transcriptome of hPSCs-derived progenitors to their in vivo counterparts.
73 To address these issues, we generated scRNAseq data sets of human hPSCs-derived
74 haematopoietic progenitors. Lineage trajectories predicted in silico were validated by
75 functional assays of sorted cells and the specificity of our marker repertoire was confirmed
76 using a CITE-seq¹⁷ strategy. Furthermore, to annotate the hPSCs-derived progeny in an
77 unbiased manner, we employed machine learning and trained an artificial neural network
78 (ANN) to recognize the single-cell gene expression profiles of human fetal liver cells. This
79 trained ANN was subsequently used to predict the identities of hPSCs-derived cells. The
80 ANN thereby provides a mapping between in vivo and in vitro hematopoiesis and revealed
81 a subset of hPSCs-derived cells that closely resembles HSCs in the foetal liver. Finally, by
82 comparing that transcriptome of in vitro and in vivo-generated HSCs we identified genetic
83 pathways that can be exploited to improve HSCs production in vitro from hPSCs.

84

85 **Methods**

86 Methods are available as supplementary methods.

87

88 **Results**

89

90 **Single cell RNA sequencing of iPSCs-derived haematopoietic progenitor cells** 91 **reveals the transcriptome of naïve and lineage committed progenitors.**

92 To resolve the heterogeneity of in vitro generated hPSCs-derived haematopoietic
93 progenitors we designed a minimal membrane marker strategy that allows to broadly
94 isolate hPSCs-derived haematopoietic progenitors. This marker strategy was validated
95 using two hPSC reporter lines, RUNX1C-GFP and KLF1-mCherry, CFU-C assays of
96 sorted cell populations and gene expression profiling (Supplementary Figure S1). These
97 data, together with previous reports^{18–22}, supported our rationale that the isolation of
98 CD235a⁺CD43⁺ cells would enrich for HSPCs and exclude cells derived from the primitive
99 wave (Supplementary Figure S1). We anticipated that the CD235a⁺CD43⁺ compartment

100 would also comprise the early stages of lineage commitment, capturing the downstream
101 hierarchy of early human progenitors.
102 CD235a-CD43⁺ suspension cells from two independent replicate cultures at day 13 of
103 differentiation were sorted by FACS and subjected to scRNAseq and data analyses
104 (Figure 1A). After quality control and clustering²³ we obtained the transcriptome of 11420
105 cells (Supplementary Figure 2A-C) belonging to **eight** clusters of cells (Figure 1B).
106 Although the two replicates did not show obvious differences (Supplementary Figure 2C),
107 any potential batch effect was regressed out prior to integration. We assigned cell
108 identities based on the expression of known markers and identified additional markers
109 from the dataset that were cluster specific (Figure 1C-D). Clusters containing more
110 immature, uncommitted progenitors were identified by their expression of progenitor-
111 associated genes such as *KIT* and *GATA2* and their lack of expression of genes
112 associated with specific cell lineages and were thus annotated as naïve populations
113 (Figure 1D, Supplementary Figure 2D). Clusters that displayed expression of lineage
114 markers were annotated as lineage committed progenitors (Figure 1B-D, Supplementary
115 Figure 2D), including clusters of cells committed towards the megakaryocyte (*GP9* and
116 *PF4*), erythroid (*GYP A* and *KLF1*) and granulocyte (*AZU1* and *PRNT3*) lineages (Figure
117 1D). Markers for each of the cell clusters were identified by differential gene expression
118 analysis, further supporting the identities assigned to each cluster (Figure 1C,
119 Supplementary Table 1).

120

121 **Trajectory analyses reveal the hierarchy of in vitro derived haematopoietic** 122 **progenitors.**

123 To study the hierarchical relationship between cell populations, we performed trajectory
124 analysis using different methods including diffusion analysis²⁴ using Seurat R package²³
125 and pseudotemporal ordering, using Monocle R package²⁵ and Partition-based graph
126 abstraction (PAGA)²⁶. Diffusion analysis identified a central core from which three distinct
127 trajectories emerged (Figure 2A). The central core corresponded to cells that we had
128 annotated as naïve progenitors (Figure 2A-B). Branches comprised cells expressing genes
129 associated with specific lineages, annotated as Erythroid (Ery)-, Megakaryocyte (Mega)-
130 and Granulocyte (Granulo)-committed lineages. Comparable trajectories were observed
131 using pseudotemporal ordering with PAGA and Monocle (Figure 2C-D). The PAGA
132 analysis showed that the naïve cells were highly connected to erythroid, megakaryocyte

133 and granulocyte committed cells (Figure 2C). Pseudotime reconstruction of the hierarchy
134 showed that cells annotated as naïve 1 were located at the top of the hierarchy and
135 appeared to progress to naïve 2 cells before entering branches containing lineage
136 committed cells (Figure 2D-E). Lineage commitment was also inferred from the expression
137 of lineage-associated transcription factors that were filtered from the cluster specific
138 marker genes according to their GO annotation (Figure 2F, **Supplementary Figure 2E**). For
139 example, erythroid committed clusters demonstrated expression of both *KLF1* and *MYC*,
140 with the latter decreasing in Ery 2 compared to Ery1, in keeping with their position within
141 the hierarchy (Figure 2E, F). Within the megakaryocyte-committed clusters 1 and 2 we
142 observed the expression of *GATA1*, *TAL1* and *FLI1* a cocktail of genes recently used for
143 hPSCs forward programming to megakaryocytes (Figure 2F)²⁷. Granulocyte-committed
144 cells were represented by a separate branch and demonstrated the expression of *CEBP-*
145 *D*, *CEBP-B*, *CEBP-A* and *CEBP-E* (Figure 2F). We then focused our attention on the
146 transcription factors expressed by the naïve progenitor clusters and noted a high level of
147 expression of *LMO4* and *ID2*, as well as *GATA2* which is known to be expressed in
148 HSPCs (Figure 2F). We then inferred their role in the gene network using a network
149 inference algorithm (Partial Information Decomposition and Context, PIDC)^{28,29}. Single cell
150 transcriptomic data offers the potential to observe dependencies between the expression
151 profiles of pairs of genes, that if co-regulated, are expected to change in a coordinated
152 fashion. Genes with highest statistical dependencies are connected by edges that
153 altogether define the network^{28,29}. Many of the transcription factors previously detected to
154 be highly expressed within the naïve cell populations such as *ID2*, *ID4* and *LMO4*, occupy
155 nodes within this large network (Supplementary Figure 3A-B). This strategy corroborates
156 the importance of the identified transcription factors as functional elements within the
157 single cell gene network.

158

159 **CD44 membrane expression marks human clonogenic haematopoietic progenitors.**

160 To experimentally validate the results of our trajectory analyses experimentally, we set out
161 to assess the haematopoietic potential of the naïve progenitor populations. We defined a
162 prospective sorting strategy using genes encoding the cell surface markers *CD33*, *CD44*,
163 and *ITGB2* (also known as CD18) that were enriched within the naïve progenitors' clusters
164 (Figure 3A). *CD33* was expressed by both naïve 1 and naïve 2 progenitors whereas *CD44*
165 and *CD18* expression appeared higher in the naïve 1 population (Figure 3A). We

166 fractionated CD235a⁻CD43⁺CD33⁺ cells and identified subpopulations as naïve 1A
167 (CD44⁺CD18⁻), naïve 1B (CD44⁺CD18⁺) and naïve 2 (CD44⁻CD18⁻) (Figure 3B). Trajectory
168 analysis predicted that naïve 1 cells were at the top of the hierarchy and gave rise to the
169 naïve 2 cells prior to lineage commitment (Figure 2D-E). To test this in silico prediction, we
170 used a chimeric co-culture system using the Zeiss Green (ZsG) reporter (Figure 3C). This
171 approach allowed us to sort, for example, ZsG-labelled naïve 1 cells, then track their ZsG
172 progeny after being placed back in the complex differentiation environment. We
173 synchronously differentiated the ZsG-iPSC line, constitutively expressing the fluorescent
174 reporter³⁰, and the parental iPSC line. To verify the progressions of naïve 1 to naïve 2 and,
175 naïve 2 to lineage committed cells, we sorted naïve 1 (CD33⁺CD44⁻CD18⁻) or naïve 2
176 (CD33⁺CD44⁺CD18^{-/+}) cells from ZsG-iPSCs at day 10 and co-cultured these with the
177 synchronized differentiating parental cells for a further 3 days. As predicted from the
178 trajectory analysis, the naïve 1 cell population was able to generate ZsG-expressing naïve
179 2 cells. We also noted that the naïve 1 cells retained their immunophenotype, indicating
180 some self-renewal capacity (Figure 3D). Interestingly, naïve 2 cells demonstrated some
181 potential to acquire CD44 and CD18, markers of naïve 1 cells (Figure 3D), suggesting
182 fluidity between these states. As predicted by our trajectory analyses (Figure 2D-E), naïve
183 2 cells acquired the ability to generate more mature cells including erythroid cells
184 (CD235a⁺), megakaryocytes (CD41⁺) and macrophages (25F9⁺) (Supplementary Figure
185 3C). We compared the colony forming capacity of naïve 1 and 2 progenitors present at day
186 10 and day 13. When plated in clonogenic CFU-C assays, CD44⁺ naïve 1 cells formed
187 CFU-C colonies but virtually no colonies were generated by naïve 2 cells at either time
188 point (Figure 3E-F). These data support the proposed hierarchy and indicate that CD44
189 expression alone resolves colony forming cells. Our chimeric co-culture system was able
190 therefore to assess the lineage output that could not be assessed solely by CFU-C assays.
191 We observed that naïve progenitors expressed high levels of ID genes (Figure 2F), and
192 that they were identified as nodes within the gene network (Supplementary Figure 3A). As
193 ID genes are targets of BMP signaling, we predicted that these naïve cells would be
194 responsive to BMP stimulation. We added BMP4 to differentiation culture from day 10,
195 when both naïve 1 and 2 were present and then assessed the proportion of these cells 3
196 days later. In presence of BMP4, we observed a 25% and 59% expansion of naïve 1 and 2
197 cells respectively (Supplementary Figure 3E). In this experiment we used both hESCs and
198 hiPSCs and showed that naïve progenitors are present at a comparable frequency in both

199 hESCs and hiPSCs (Supplementary Figure 3D), and that naïve progenitors derived from
200 both lines responded to BMP stimulation in a comparable manner (Supplementary Figure
201 3D-E). Thus this experiment not only identified an important functional signaling pathway
202 but also confirmed that the markers we used to define naïve progenitors, and their
203 response to BMP signaling, are not PSC line specific.

204 To assess whether the naïve cell populations identified using our unique sorting strategy
205 showed features of definitive haematopoietic progenitors, we assessed the expression of
206 the RUNX1C-GFP reporter. We observed RUNX1C-GFP expression in both cell types,
207 with a higher proportion of RUNX1C⁺ cells in the naïve 1 compared to naïve 2 population
208 (Figure 3G). Definitive HSPCs are generated via endothelial to hematopoietic transition
209 (EHT) during embryonic development^{31,32} so, we would expect comparable hPSCs-derived
210 progenitors to have hallmarks of their endothelial origin. Here we demonstrate that naïve
211 CD44⁺ cells generated in vitro from hPSCs co-expressed CD34 and the endothelial marker
212 CD144 (also known as VeCad) with expression being higher at day 10, when the majority
213 of naïve progenitors were present (Supplementary Figure 3L). This endothelial signature of
214 naïve progenitors, together with their lineage potential reflects their definitive features. To
215 confirm that CD44 expression was associated with HSPCs in vivo we demonstrated its co-
216 localization with CD45 and CD144 in the mouse yolk sac and AGM region (Supplementary
217 figure 3F-J). At E10.5 in the yolk sac, CD44 was expressed on endothelial cells in a
218 bimodal fashion, with vessels expressing low and high levels, the latter being associated
219 with bright clusters of haematopoietic cells (Supplementary figure 3G). By flow cytometry,
220 we observed that by E11, all CD45⁺ cells and a proportion of CD144⁺ cells were within the
221 CD44^{high} population (Supplementary figure 3F). Within the embryo proper, CD44 was
222 expressed on the membrane of endothelial cells within the dorsal aorta, whereas venous
223 endothelial layers were CD44 negative (Supplementary figure 3H-I). CD44 was also co-
224 expressed with CD45⁺ in the AGM region (Supplementary figure 3H-J). Furthermore,
225 expression of LMO4 in CD44⁺ cells within the AGM region is in keeping with its high level
226 of expression in naïve hPSCs-derived HSPCs (Figure 2F) and supports its identification as
227 a novel haematopoietic transcription factor. These data suggest that CD44 is expressed
228 on haemogenic endothelial cells and it is retained on emerging haematopoietic cells in
229 vivo, similar to what we have observed during the in vitro differentiation of human
230 progenitors (Supplementary Figure 3I).

231

232 Identification of membrane markers of lineage committed progenitors

233 We next selected membrane markers that we predicted could be used for the isolation of
234 lineage committed progenitors. Erythroid-primed clusters 1 and 2 both showed expression
235 of *MYC* (Figure 2F) and *EPCAM* (also known as CD326) (Supplementary Figure 4A),
236 indicative of early committed erythroid cells^{33,34}. We confirmed that CD326 was expressed
237 in the majority of CD235a⁺ cells at day 13 of iPSC differentiation but interestingly, we
238 noted a small number of CD326⁺CD235a⁻ (Supplementary Figure 4B), suggesting that
239 CD326 might be marking commitment to the erythroid lineage prior to CD235a acquisition.
240 To test this, we assessed the expression dynamics of these markers during the erythroid
241 differentiation of umbilical cord blood CD34⁺ (UCB34⁺) cells. At day 10 of differentiation,
242 CD326 was expressed in CD235a^{-/low} cells but not in CD235a^{high} cells, the latter
243 corresponding to more mature erythroid cells (Supplementary Figure 4B). CD326 was not
244 expressed in cells at day 18 of the differentiation protocol (when the majority of cells are
245 mature CD235a⁺ cells) nor in the mature erythrocytes found in adult peripheral blood
246 (Supplementary Figure 4B). Taken together these data suggest that CD326 marks early
247 erythroid progenitors in both hiPSC, fetal and adult derived cells. We also noted the
248 expression of *HBG1*, *HBG2*, *HBA1*, and *HBA2*, subunits of fetal hemoglobin, indicative of
249 erythroid cells derived from definitive hematopoiesis (Supplementary Figure 4C).

250 Three clusters with megakaryocyte and platelet signatures (Mega-primed 1, 2 and 3) were
251 predicted by expression of *ITGA2B* (CD41), *GP9*, *PF4* (Figure 1C-D and Supplementary
252 Table 1). *ICAM2* was expressed at higher level in cluster Mega-primed 3 (Supplementary
253 Figure 4D), as for CD9, known to increase along megakaryocytes differentiation³⁵. *ICAM2*
254 and CD9 co-expression was confirmed by flow cytometry (Supplementary Figure 4D). We
255 observed a population of CD41⁺CD9⁺*ICAM2*⁺ cells, with around 85% of the CD41⁺CD42a⁺
256 (Supplementary Figure 4E), that did not detect polyploidy, supporting their immature status
257 (Supplementary Figure 4F-G).

258 Granulocyte-committed clusters were predicted by expression of markers such as *MPO*,
259 *AZU1*, *RNASE2* and *ITGB2* which encodes the membrane marker CD18, subunit of
260 different leukocytes marker such as CD11a-d, Mac-1 and LFA-1 (Figure 1C,
261 Supplementary Table 1). We sorted CD235a⁻CD43⁺CD33⁺CD44⁻CD18⁺ cells and
262 confirmed the phenotype of granulocytes and monocytes based on their nuclear
263 morphology (Supplementary Figure 4H). Further clustering revealed three sub-clusters
264 corresponding to eosinophil, neutrophils and monocytes lineages (Supplementary Figure

265 4I-L). Noteworthy, *RUNX3* expression was specifically associated with the monocyte
266 subcluster (Supplementary Figure 4J) previously reported to be important for zebrafish
267 stem cells and macrophages³⁶, and to be expressed by mouse tissue resident
268 macrophages of the skin³⁷.

269 In summary, we showed that naïve progenitors give rise also to committed progenitors
270 characterized by features of granulocytes and monocyte, cell types that emerge
271 exclusively in the definitive waves⁵.

272

273 **CITE-seq analysis of human iPSC-derived haematopoietic progenitors.**

274 To further study the temporal emergence of the progenitor populations and their
275 associated markers, we carried out CITE-seq analysis whereby single cell membrane
276 marker expression can be directly correlated with the relative transcriptome¹⁷. To ensure
277 that we sampled even the rarest progenitor cell population we extended the CITE-seq
278 analysis to adherent cells and included an earlier time point (day 10) in addition to day 13.
279 Again, to exclude primitive erythroid cells, we selected CD235a-negative suspension cells
280 but, in this experiment, we included and enriched for CD43⁻ cells that had been excluded
281 in our initial study. (Supplementary figure 5B). **We expected early progenitors to express**
282 **CD31 and to potentially remain part of hematopoietic clusters within the adherent fraction**
283 **of the culture and so we FAC-sorted the adherent cells into CD31⁻ and CD31⁺ fractions.**

284 Cells were labeled with oligonucleotide tagged antibody specific for the membrane
285 markers that we identified in our initial experiment (ADT_CD18, ADT_CD33, ADT_CD41,
286 ADT_CD44, ADT_CD102, ADT_CD326; **ADT: Antibody-Derived Tag**) as well as other
287 markers of endothelial and early committed hematopoietic cells (ADT_CD144) and of
288 macrophages (ADT_CD163). To test the specificity of the membrane marker repertoire
289 previously identified on the suspension cells, we subset the two libraries corresponding to
290 suspension cells collected at day 10 and 13 (Figure 4, Supplementary Figure 5B-C). After
291 multidimension reduction and clustering analysis we identified a naïve progenitor
292 population (Figure 4A), comparable to our first sequencing experiment (Figure 2A). These
293 naïve progenitors exhibited erythroid (Ery), megakaryocyte (Mega), and granulocyte and
294 monocytes (Gra-Mo) lineage potential, with increased lineage commitment at day 13
295 compared to day 10 (Figure 4B); in line with the expression pattern of genes associated
296 with naïve and committed stages in these days (Supplementary Figure 5D). Analysis of the
297 ADTs showed that each marker was expressed in the expected cluster (Figure 4C) thus

298 supporting them as markers for defined progenitors. To further explore the power of the
299 ADT approach, we performed multidimension reduction using ADTs as the only input
300 dimensions and proved that ADT data alone identified remarkably similar clusters (Figure
301 4D-E), strongly correlated with the clusters derived from the entire transcriptome (Figure
302 4F). Taken together, the CITE-seq approach confirms that the markers identified from our
303 first scRNAseq analysis define the hierarchy of human developmental hematopoiesis in
304 vitro with high specificity.

305

306 **Comparison of in vitro generated progenitors with in vivo produced cells.**

307 **The use of human PSCs as a renewable source of hematopoietic cell types faces major**
308 **challenges relating to, for example, the inefficient repopulation capacity of progenitor cells**
309 **and the incomplete maturation of differentiated cell types. To identify the underlying**
310 **molecular basis associated with these deficiencies in hPSC-derived cells, we compared**
311 **our dataset to a human fetal liver dataset which contains the complete hematopoietic**
312 **hierarchy from long-term reconstituting HSCs to differentiated cell types.**

313 To assess how hPSCs-derived naïve and lineage-committed progenitors compared to their
314 equivalent counterpart generated in vivo, we assessed the expression of selected genes
315 identified to distinguish the various cell types detected in the human fetal liver³⁸ (Figure
316 4G). An initial analysis of marker genes of lineage commitment in the developing embryo
317 revealed that these markers are remarkably powerful for discriminating the equivalent in
318 vitro cell types identified in our in vitro study (Figure 4G, Supplementary Table 1).

319 Interestingly, *SPINK2*, a newly reported marker of fetal HSC/MPP³⁸, was also expressed
320 specifically by our naïve progenitor cells (Figure 4G), together with CD34 (Supplementary
321 figure 3L). **These specific similarities observed between in vitro and in vivo developing**
322 **hematopoietic progenitor cells led us to investigate in a more comprehensive manner the**
323 **phenotype of cell types that are produced in vitro and how well these in vitro derived cells**
324 **reflect the corresponding cell types during in vivo development. Therefore, we used the**
325 **same published human fetal liver scRNAseq data as a reference, firstly, to identify in vitro**
326 **derived cells with gene expression signatures of human fetal liver hematopoietic cells and,**
327 **secondly, to quantify the similarity to their corresponding transcriptomes.** To address the
328 first question, we employed machine learning to transfer labels from the fetal liver
329 reference data to our in vitro-derived blood cells (Figure 5A). This approach enabled a
330 much broader and unbiased identification of cell types compared to inference based purely

331 on marker genes. We followed our recently developed strategy³⁹ and trained an artificial
332 neural network (ANN)³⁹ to recognize single-cell gene expression profiles of human foetal
333 liver cells that were sampled at a time in development at which the liver is the main site of
334 blood cell formation³⁸. Briefly, this ANN is trained using the expression data of 3,479
335 genes and 145,725 cells from fetal liver as an input³⁸. From these labelled data, the ANN
336 learns to predict, from which of the 28 different fetal liver cell types a particular gene
337 expression pattern originates. Once trained, the ANN is given previously unseen test data
338 from in vitro derived cells as an input in order to annotate these data with human fetal liver
339 cell labels. Since this approach considers 3,479 genes, it enabled a more comprehensive
340 identification of cell types based on similarities in global gene expression patterns rather
341 than specific marker genes.

342 The ANN was able to identify cell types within the source domain (the fetal liver data) with
343 high accuracy as shown by the performance metrics obtained from 5-fold cross-validation
344 (Supplementary figure 6A-B). The trained ANN was subsequently applied to the target
345 domain (in vitro) to test if the hPSCs-derived cells were similar to those present in the
346 foetal liver, in which case the label of that specific in vivo cell would be transferred. The
347 ANN was able to assign labels to 92% of in vitro produced cells into various cell types
348 present in vivo (Supplementary figure 6 C-D), most notably, a small population was
349 labeled as HSC/MPP. This indicates that the global gene expression pattern of a subset of
350 the in vitro derived cells is very similar to HSC/MPPs from the in vivo reference data in
351 fetal liver. To quantify precisely how similar these in vitro derived HSC/MPPs are to their in
352 vivo counterparts, we calculated the average pairwise Euclidean distance between
353 HSC/MPPs, using the human fetal liver as a reference. This analysis indicates that fetal
354 liver HSC/MPPs are, on average, only marginally more similar to one another as they are
355 to iPSC derived HSC/MPPs (Supplementary Figure 8A). In summary, this analysis
356 indicates that the in vitro derived HSC/MPPs closely, yet not perfectly, reflect the gene
357 expression patterns of their in vivo counterparts. Using the ANN we also observed that the
358 relative abundance of the predicted HSC/MPP population decreased with time by day 13
359 (Figure 5B), whereas, the relative abundance of committed cells increased over this time
360 as expected (Supplementary Figure 6E). When we applied the same ANN strategy to our
361 first data set, that was generated from day 13 progenitors that were selected on the basis
362 of CD43 expression, no HSC/MPP were detected (Figure 5C). This is consistent with our
363 observation that this transient HSC/MPP population is present in higher numbers earlier at

364 day 10, when they are almost equally distributed in the adherent CD31+ and suspension
365 CD235a- compartment (Supplementary Figure 6F). We looked for marker genes that
366 defines this predicted HSC/MPP cell population in vitro and looked specifically for
367 membrane markers according to their GO annotation (Supplementary Table 1). Together
368 with expected markers such as *CD34*, *CD44* and *CD33*, we also detected *CD132*, *CD52*,
369 *CD180* and *IL3RA* and many others that will allow to design a prospective sorting strategy
370 to isolate this specific population. We then subset the in vivo and in vitro HSC/MPP and
371 integrated the two datasets (Figure 5D). The integrated data allowed for direct comparison
372 of their transcriptome and identified 54 differentially expressed genes (Supplementary
373 Table 1), all of which were lower in HSC/MPP produced in vitro compared to those
374 generated in vivo. GO analysis of these genes identified enrichment for KEGG signaling
375 pathways such as NOD-like receptor, IL-17, NF-Kappa B and HIF-1 (Supplementary Table
376 1). We also identified 6 genes encoding transcription factors: *EGR1*, *ZFP36L1*, *NR4A1*,
377 *FOS*, *JUN* and *JUNB* (Figure 5E). Interestingly, the EGR1 binding site was enriched,
378 amongst others, in the upstream region of the differentially expressed genes (Figure 5F),
379 suggesting an important regulatory role of EGR1.

380 We also compared the predicted HSC/MPP derived from hPSCs to hematopoietic
381 progenitors isolated from different sites of hematopoiesis in the developing embryo
382 including to fetal liver HSC/MPP³⁸, yolk sac MPP³⁸ that were collected at Carnegie stages
383 5 to 14, and AGM⁴⁰ sorted progenitors (CD34+CD45+CD235a-) collected at Carnegie
384 stage 15, around the time of early HSC emergence (Supplementary Figure 7A-B). Whole
385 transcriptome comparison, followed by KEGG pathway analysis, showed that in vitro
386 HSC/MPP cells are marked by genes associated with oxidative phosphorylation
387 (Supplementary Table 1), indicating metabolic differences between in vitro and in vivo
388 produced progenitors. Hypoxic conditions characterize mammalian embryo
389 development⁴¹, and more specifically the development of the hematopoietic system,
390 where hypoxia has been detected in hematopoietic clusters in the AGM region, and in the
391 fetal liver⁴². The hematopoietic progenitors derived from hPSCs were instead
392 differentiated in normoxic conditions which could explain their different metabolic profile.
393 The fetal liver cells were marked by HLA genes and consequently KEGG pathways
394 associated with antigen presentation and T-cell development (Supplementary Table 1).
395 The AGM dataset displayed high expression levels of genes associated with Notch
396 pathway, such as *HES1*, *NOTCH1*, *NOTCH2*, *JAG1* and *JAG2*. This is in line with the

397 developmental stage at which they were collected when the Notch pathway is
398 orchestrating the HSC emerge⁴³. Within the markers of yolk sac progenitors, we detected
399 genes related to early hematopoietic development. *FRZB*, mesodermal cell marker, and
400 *HBE1*, marker of primitive hematopoiesis, were listed in the top 10 differentially expressed
401 genes: this underlines the early developmental features of yolk sac progenitors. Finally, we
402 noted also that *SPINK1* was identified as marker for YS progenitors. While *SPINK2*,
403 identified here and by others as marker of progenitor cells^{38,40} was expressed by
404 progenitors from all the tissues, *SPINK1* was detected exclusively in the YS
405 (Supplementary Figure 7X), suggesting that this gene could discriminate extraembryonic
406 from intraembryonic hematopoiesis.

407 Finally, we compared lineage committed cells identified by the ANN, in our in vitro dataset,
408 to their in vivo counterpart from fetal liver, to identify genes that can be used as targets to
409 improve the production in vitro of differentiated blood cell types. We listed the differentially
410 expressed genes between in vitro and fetal liver cells and identified the transcription
411 factors within the list (Supplementary Figure 8, Supplementary Table 1). Particularly
412 interesting, late erythroid cells in vitro show high level of *PCLG2*, phospholipase C gamma
413 2, able to control intracellular calcium via production of IP3, inositol triphosphate.
414 Intracellular calcium peaks just before enucleation, prior nuclei extrusion in the
415 orthochromatic erythroblasts⁴⁴. Erythroid cells derived from hPSCs are characterized by a
416 general inefficient enucleation^{45,46}, independently or their primitive or definitive origin
417 (Supplementary Figure 1J) and this could be related to their intracellular calcium control.

418 In summary, we have identified a rare population of HSC/MPP-like cells in vitro that
419 emerge early during differentiation of hPSCs and that display broadly similar gene
420 expression patterns when compared to HSCs in development. However subtle differences
421 are also apparent and a more detailed study of these differences could explain the known
422 deficiencies of PSC-derived cells and ultimately be exploited to improve their therapeutic
423 use. Our novel approach combines scRNAseq and machine learning to help identify
424 candidate genes that may improve the production of HSCs and mature lineage cells from
425 pluripotent stem cells in vitro, by closely recapitulating in vivo hematopoiesis.

426

427

428

429 **Discussion**

430 We described the single cell transcriptome and membrane markers of naïve hematopoietic
431 progenitors and their lineage committed descendants derived in vitro from human
432 pluripotent stem cells. The repertoire of membrane markers proved to be remarkably
433 accurate in capturing the different states prior to and after lineage commitment.
434 We identified a population of naïve progenitors situated at the top of the differentiation
435 hierarchy, marked by CD44, a protein involved in the hematopoietic transition of the
436 hemogenic endothelium in the mouse AGM region⁴⁷. We validated their lineage potential
437 employing a chimeric culture system, where isolated naïve progenitors, marked by Zeiss-
438 Green expression, demonstrated overlapping lineage output to that predicted in silico.
439 We also observed that progenitors are capable of moving between the naïve states, as
440 well as progressing into committed states. This is in keeping with many other scRNAseq
441 and proteomic studies that have reported a continuum of cell states as opposed to
442 sequential discrete cell types hierarchies⁴⁸⁻⁵¹. In line with a recent murine study⁴⁷, we have
443 shown that CD44 is expressed in naïve hPSCs-derived progenitors and here we
444 demonstrated that both human and mouse progenitors also express LMO4, a LIM-domain
445 protein⁵². Recent scRNAseq detected LMO4 in both human granulocyte progenitors in the
446 bone marrow⁴⁸ and adult mouse HSC⁵³, but its associated proteins have not been
447 identified. We also reported high levels of *ID* genes within the progenitors, target genes of
448 BMP signaling known to be involved in HSC emergence⁵⁴⁻⁵⁶. IDs, like LMOs proteins, do
449 not present DNA binding domain and rather act through binding of other proteins in
450 complexes also involved in HSPC development⁵⁷ and erythropoiesis⁵⁸. Overexpression of
451 ID2 in human HSC from cord blood has been reported to enhance their functional
452 stemness in vivo⁵⁹, supporting the idea that this class of proteins might maintain the
453 progenitor status and thus might be useful in alternative programming strategies of hPSCs.
454 The use of scRNAseq on vast numbers of cells allows to detect even the rarest cell
455 population and we considered that it might enable the detection of rare HSC-like cells in
456 differentiating hPSCs cultures. We showed to hPSCs-derived cells showed a remarkable
457 specific expression pattern of marker genes identified in the human embryo, for example,
458 *SPINK2*, a novel marker of human fetal liver HSC and MPP, marked also our naïve
459 progenitors. By using machine learning we identify specific cell types sampled in vivo and
460 detected a small and transient population of HSC-like cells that, when compared to their in
461 vivo counterpart from fetal liver, showed only small transcriptional differences. Previous
462 reports described the hematopoietic progenitors obtained with the differentiation employed

463 in this work as intraembryonic-like¹⁸, using T-cells lineage as hallmark of definitive
464 hematopoiesis. However, yolk sac shows T-cell potential prior to HSC emergence^{60,61},
465 thus limiting the use of T-cell assay alone as discriminative of the corresponding
466 developmental wave. Our machine learning approach and the detection of HSC-like cells
467 strongly supports the intraembryonic identities of the hematopoietic cells differentiated in
468 vitro and provide an alternative and multifactorial approach to address questions regarding
469 the similarities to developmental counterparts. The unbiased and comprehensive
470 comparison used in this study allowed us to pinpoint differentially expressed genes
471 between in vitro-derived and in vivo HSCs can now be exploited to improve production of
472 HSC in vitro. Our analysis indicates in vitro HSC-like cells do not express CD43,
473 comparable to mouse Pro-HSC prior their maturation into functional definitive HSC⁶². This
474 could suggest that the widely acknowledged inability of hPSCs-derived progenitors to
475 reconstitute the hematopoietic system, could be due to their immature phenotype and the
476 lack of appropriate culture conditions for HSCs maturation and maintenance. In addition,
477 the identification of the HSC-like population using our machine learning approach, which is
478 based on high similarity in the gene expression profiles, could suggest that the molecular
479 basis of the functional deficiency of this in vitro derived population could reside at a post-
480 transcriptional level. Thus, future experiment will be required to assess whether a further
481 ad-hoc maturation step of sorted HSC-like cells would achieve reconstitution.

482 When we compared the hematopoietic progenitors developed in the human embryo
483 throughout gestations together with those derived in vitro, we found that while *SPINK2*
484 was expressed by all progenitors, *SPINK1* was exclusively detected in cells from the yolk
485 sac. *SPINK1* is able to bind to EGFR and induce epithelial to mesenchymal transition in
486 cancer cells^{63,64}, a process similar to the endothelial to hematopoietic transition, where the
487 role of *SPINK1* remains largely unexplored. In summary, we propose here *SPINK1* as a
488 possible marker for primitive hematopoiesis which could be an extremely useful genes to
489 trace the cells that colonize the embryo from the yolk sac.

490 The differentiation protocol used in this study is well defined and serum-free and is one of
491 the most commonly used protocols used by many laboratories. We also showed that our
492 markers are able to identify functionally similar progenitors in different cell lines. Thus, our
493 browsable datasets and the findings of our study will be of interest to many in the field of
494 hematopoiesis and will allow to test how the frequency of this populations vary in response
495 to different cytokines conditions. In addition, the increasing availability of large scRNAseq

496 dataset of human tissue makes our pipeline applicable to the analyses of other systems
497 where the hPSCs differentiation aims to produce adult-like cells for clinical application. In
498 this way cell types differentiated in vitro can now be annotated in an unbiased manner that
499 does not rely on few known markers and allows the identification of transcriptional
500 discrepancies between cell types produced in vitro and their in vivo counterparts.
501 In conclusion, our browsable dataset provides a comprehensive transcriptional
502 characterization of in vitro derived hematopoietic progenitors. This work defines the
503 makeup of the progenitor populations that give rise to immune cells, such as macrophages
504 and granulocytes, as well as HSC-like cells, which holds great promise for their
505 therapeutically application.

506

507 **Acknowledgment**

508 The work was funded by Wellcome Trust (Grant No. 102610), MRC Innovate UK (Grant
509 No. 102853), BBSRC (Grant No. S002219/1). AF received a Carnegie Incentive Grant
510 (Grant No. RIG008218). AB received BBSRC Future Leaders Fellowship (Grant reference
511 BB/N011597/1). Sequencing was carried out by Edinburgh Genomics, The University of
512 Edinburgh. Edinburgh Genomics is partly supported through core grants from NERC
513 (R8/H10/56), MRC (MR/K001744/1) and BBSRC (BB/J004243/1). We thank Professor
514 Ben Macarthur for suggesting the application of machine learning to our study; Andrew
515 Elefanty for sharing the RUNX1C-GFP cell line; Fiona Rossi, Claire Cryer, Bindi Heer and
516 Andrea Corsinotti from the Flow Facility as well as Bertand Verney and Matthieu Vermeren
517 from the imaging facility. This work has made use of the resources provided by the
518 Edinburgh Compute and Data Facility (ECDF) (<http://www.ecdf.ed.ac.uk/>).

519

520 **Author Contribution**

521 AF, designed and performed research, analyzed the data and wrote the manuscript. AF,
522 PS, AB and NR performed bioinformatics analysis. PR, ST, MLY, AHT, JE, BH, RA
523 performed research. LMF designed the experiment, analyzed data and wrote the
524 manuscript. NH, AM, KO, and NR provided intellectual input and final approval of the
525 manuscript.

526

527 **Declaration of interest**

528 Authors declare no competing interests.

529

530 **References**

531

- 532 1. Vo LT, Daley GQ. De novo generation of HSCs from somatic and pluripotent stem
533 cell sources. *Blood*. 2015;125(17):2641–8.
- 534 2. Palis J. Hematopoietic stem cell-independent hematopoiesis: emergence of
535 erythroid, megakaryocyte, and myeloid potential in the mammalian embryo. *FEBS*
536 *Lett*. 2016;590(22):3965–3974.
- 537 3. Tober J, Koniski A, McGrath KE, et al. The megakaryocyte lineage originates from
538 hemangioblast precursors and is an integral component both of primitive and of
539 definitive hematopoiesis. *Blood*. 2007;109(4):1433–41.
- 540 4. Palis J, Robertson S, Kennedy M, Wall C, Keller G. Development of erythroid and
541 myeloid progenitors in the yolk sac and embryo proper of the mouse. *Development*.
542 1999;126(22):5073–84.
- 543 5. McGrath KE, Frame JM, Fegan KH, et al. Distinct Sources of Hematopoietic
544 Progenitors Emerge before HSCs and Provide Functional Blood Cells in the
545 Mammalian Embryo. *Cell Rep*. 2015;
- 546 6. Mass E, Ballesteros I, Farlik M, et al. Specification of tissue-resident macrophages
547 during organogenesis. *Science (80-.)*. 2016;353:6304.
- 548 7. Schulz C, Gomez Perdiguero E, Chorro L, et al. A lineage of myeloid cells
549 independent of Myb and hematopoietic stem cells. *Science (80-.)*.
550 2012;336(6077):86–90.
- 551 8. Stremmel C, Schuchert R, Wagner F, et al. Yolk sac macrophage progenitors traffic
552 to the embryo during defined stages of development. *Nat. Commun*. 2018;9(1):75.
- 553 9. Medvinsky A, Dzierzak E. Definitive hematopoiesis is autonomously initiated by the
554 AGM region. *Cell*. 1996;86(6):897–906.
- 555 10. Ivanovs A, Rybtsov S, Welch L, et al. Highly potent human hematopoietic stem cells
556 first emerge in the intraembryonic aorta-gonad-mesonephros region. *J. Exp. Med*.
557 2011;208(12):2417–2427.
- 558 11. Easterbrook J, Fidanza A, Forrester LM. Concise review: Programming human
559 pluripotent stem cells into blood. *Br. J. Haematol*. 2016;173(5):.
- 560 12. Ivanovs A, Rybtsov S, Anderson RA, Turner ML, Medvinsky A. Identification of the
561 niche and phenotype of the first human hematopoietic stem cells. *Stem Cell Reports*.
562 2014;2(4):449–456.
- 563 13. Tavian M, Hallais MF, Péault B. Emergence of intraembryonic hematopoietic
564 precursors in the pre-liver human embryo. *Development*. 1999;126(4):793–803.
- 565 14. Guibentif C, Rönn RE, Böiers C, et al. Single-Cell Analysis Identifies Distinct Stages
566 of Human Endothelial-to-Hematopoietic Transition. *Cell Rep*. 2017;19(1):10–19.
- 567 15. Angelos MG, Abrahante JE, Blum RH, Kaufman DS. Single Cell Resolution of
568 Human Hematoendothelial Cells Defines Transcriptional Signatures of Hemogenic
569 Endothelium. *Stem Cells*. 2018;36(2):206–217.
- 570 16. Han X, Chen H, Huang D, et al. Mapping human pluripotent stem cell differentiation
571 pathways using high throughput single-cell RNA-sequencing. *Genome Biol*.
572 2018;19(1):1–19.
- 573 17. Stoeckius M, Hafemeister C, Stephenson W, et al. Simultaneous epitope and
574 transcriptome measurement in single cells. *Nat. Methods*. 2017;14(9):865–868.
- 575 18. Sturgeon CM, Ditadi A, Awong G, Kennedy M, Keller G. Wnt signaling controls the
576 specification of definitive and primitive hematopoiesis from human pluripotent stem

- 577 cells. *Nat. Biotechnol.* 2014;32(6):554–561.
- 578 19. Vodyanik MA, Thomson JA, Slukvin II, Dulac C, Péault B. Leukosialin (CD43)
579 defines hematopoietic progenitors in human embryonic stem cell differentiation
580 cultures. *Blood.* 2006;108(6):2095–105.
- 581 20. Garcia-Alegria E, Menegatti S, Fadlullah MZH, et al. Early Human Hemogenic
582 Endothelium Generates Primitive and Definitive Hematopoiesis In Vitro. *Stem Cell*
583 *Reports.* 2018;11(5):1061–1074.
- 584 21. Ng ES, Azzola L, Bruveris FF, et al. Differentiation of human embryonic stem cells to
585 HOXA+ hemogenic vasculature that resembles the aorta-gonad-mesonephros. *Nat.*
586 *Biotechnol.* 2016;34(11):1168–1179.
- 587 22. Sroczynska P, Lancrin C, Kouskoff V, Lacaud G. The differential activities of Runx1
588 promoters define milestones during embryonic hematopoiesis. *Blood.*
589 2009;114(26):5279–89.
- 590 23. Butler A, Hoffman P, Smibert P, Papalexi E, Satija R. Integrating single-cell
591 transcriptomic data across different conditions, technologies, and species. *Nat.*
592 *Biotechnol.* 2018;36(5):411–420.
- 593 24. Haghverdi L, Büttner M, Wolf FA, Buettner F, Theis FJ. Diffusion pseudotime
594 robustly reconstructs lineage branching. *Nat. Methods.* 2016;13(10):845–848.
- 595 25. Qiu X, Mao Q, Tang Y, et al. Reversed graph embedding resolves complex single-
596 cell trajectories. *Nat. Methods.* 2017;14(10):979–982.
- 597 26. Wolf FA, Hamey FK, Plass M, et al. PAGA: graph abstraction reconciles clustering
598 with trajectory inference through a topology preserving map of single cells. *Genome*
599 *Biol.* 2019;20(1):1–9.
- 600 27. Moreau T, Evans AL, Vasquez L, et al. Large-scale production of megakaryocytes
601 from human pluripotent stem cells by chemically defined forward programming. *Nat.*
602 *Commun.* 2016;7(1):11208.
- 603 28. Stumpf PS, Smith RCG, Lenz M, et al. Stem Cell Differentiation as a Non-Markov
604 Stochastic Process. *Cell Syst.* 2017;5(3):268-282.e7.
- 605 29. Chan TE, Stumpf MPH, Babbitt AC. Gene Regulatory Network Inference from Single-
606 Cell Data Using Multivariate Information Measures. *Cell Syst.* 2017;5(3):251-267.e3.
- 607 30. Lopez-Yrigoyen M, Fidanza A, Cassetta L, et al. A human iPSC line capable of
608 differentiating into functional macrophages expressing ZsGreen: A tool for the study
609 and in vivo tracking of therapeutic cells. *Philos. Trans. R. Soc. B Biol. Sci.*
610 2018;373(1750):.
- 611 31. Jaffredo T, Gautier R, Eichmann A, Dieterlen-Lièvre F. Intraaortic hemopoietic cells
612 are derived from endothelial cells during ontogeny. *Development.*
613 1998;125(22):4575–83.
- 614 32. Zovein AC, Hofmann JJ, Lynch M, et al. Fate Tracing Reveals the Endothelial Origin
615 of Hematopoietic Stem Cells. *Cell Stem Cell.* 2008;3(6):625–636.
- 616 33. Jayapal SR, Lee KL, Ji P, et al. Down-regulation of Myc is essential for terminal
617 erythroid maturation. *J. Biol. Chem.* 2010;285(51):40252–65.
- 618 34. Lammers R, Giesert C, Grünebach F, et al. Monoclonal antibody 9C4 recognizes
619 epithelial cellular adhesion molecule, a cell surface antigen expressed in early steps
620 of erythropoiesis. *Exp. Hematol.* 2002;30(6):537–545.
- 621 35. Clay D, Rubinstein E, Mishal Z, et al. CD9 and megakaryocyte differentiation. *Blood.*
622 2001;97(7):1982–1989.
- 623 36. Kalev-Zylinska ML, Horsfield JA, Flores MVC, et al. Runx3 is required for
624 hematopoietic development in zebrafish. *Dev. Dyn.* 2003;228(3):323–336.
- 625 37. Fainaru O, Woolf E, Lotem J, et al. Runx3 regulates mouse TGF- β -mediated

- 626 dendritic cell function and its absence results in airway inflammation. *EMBO J.*
 627 2004;23:969–979.
- 628 38. Popescu DM, Botting RA, Stephenson E, et al. Decoding human fetal liver
 629 haematopoiesis. *Nature.* 2019;574(7778):365–371.
- 630 39. Stumpf PS, Du D, Imanishi H, et al. Mapping biology from mouse to man using
 631 transfer learning. *bioRxiv.* 2019;
- 632 40. Zeng Y, He J, Bai Z, et al. Tracing the first hematopoietic stem cell generation in
 633 human embryo by single-cell RNA sequencing. *Cell Res.* 2019;29(11):881–894.
- 634 41. Dunwoodie SL. The Role of Hypoxia in Development of the Mammalian Embryo.
 635 *Dev. Cell.* 2009;17(6):755–773.
- 636 42. Imanirad P, Solaimani Kartalaei P, Crisan M, et al. HIF1 α is a regulator of
 637 hematopoietic progenitor and stem cell development in hypoxic sites of the mouse
 638 embryo. *Stem Cell Res.* 2014;12(1):24–35.
- 639 43. Bigas A, Espinosa L. Hematopoietic stem cells: To be or Notch to be. *Blood.*
 640 2012;119(14):3226–3235.
- 641 44. Wölwer CB, Pase LB, Russell SM, Humbert PO. Calcium signaling is required for
 642 erythroid enucleation. *PLoS One.* 2016;11(1):.
- 643 45. Lopez-Yrigoyen M, Yang C-T, Fidanza A, et al. Genetic programming of
 644 macrophages generates an in vitro model for the human erythroid island niche. *Nat.*
 645 *Commun.* 2019;10(1):881.
- 646 46. Yang C-T, Ma R, Axton RA, et al. Activation of KLF1 Enhances the Differentiation
 647 and Maturation of Red Blood Cells from Human Pluripotent Stem Cells. *Stem Cells.*
 648 2017;35(4):886–897.
- 649 47. Oatley M, Bölükbaşı ÖV, Svensson V, et al. Single-cell transcriptomics identifies
 650 CD44 as a marker and regulator of endothelial to haematopoietic transition. *Nat.*
 651 *Commun.* 2020;11(1):.
- 652 48. Paul F, Arkin Y, Giladi A, et al. Transcriptional Heterogeneity and Lineage
 653 Commitment in Myeloid Progenitors. *Cell.* 2015;163(7):1663–1677.
- 654 49. Velten L, Haas SF, Raffel S, et al. Human haematopoietic stem cell lineage
 655 commitment is a continuous process. *Nat. Cell Biol.* 2017;19(4), 271–281.
- 656 50. Rodriguez-Fraticelli AE, Wolock SL, Weinreb CS, et al. Clonal analysis of lineage
 657 fate in native haematopoiesis. *Nature.* 2018;553(7687):212–216.
- 658 51. Knapp DJHF, Hammond CA, Wang F, et al. A topological view of human CD34+ cell
 659 state trajectories from integrated single-cell output and proteomic data. *Blood.*
 660 2019;133(9):927-939.
- 661 52. Grutz G, Forster A, Rabbitts TH. Identification of the LMO4 gene encoding an
 662 interaction partner of the LIM-binding protein LDB1/NLI1: a candidate for
 663 displacement by LMO proteins in T cell acute leukaemia. *Oncogene.*
 664 1998;17(21):2799–2803.
- 665 53. Lai S, Huang W, Xu Y, et al. Cell Discovery Comparative transcriptomic analysis of
 666 hematopoietic system between human and mouse by Microwell-seq. *Cell Discov.*
 667 2018;4:34.
- 668 54. Souilhol C, Gonneau C, Lendinez JG, et al. Inductive interactions mediated by
 669 interplay of asymmetric signalling underlie development of adult haematopoietic
 670 stem cells. *Nat. Commun.* 2016;(2016): 1-13.
- 671 55. Crisan M, Kartalaei PS, Vink CS, et al. BMP signalling differentially regulates distinct
 672 haematopoietic stem cell types. *Nat. Commun.* 2015;6(1):8040.
- 673 56. McGarvey AC, Rybtsov S, Souilhol C, et al. A molecular roadmap of the AGM region
 674 reveals BMPER as a novel regulator of HSC maturation. *J. Exp. Med.*

- 675 2017;214(12):3731–3751.
- 676 57. Wilson NK, Foster SD, Wang X, et al. Combinatorial Transcriptional Control In Blood
677 Stem/Progenitor Cells: Genome-wide Analysis of Ten Major Transcriptional
678 Regulators. *Cell Stem Cell*. 2010;7(4):532–544.
- 679 58. Wadman IA, Osada H, Grütz GG, et al. The LIM-only protein Lmo2 is a bridging
680 molecule assembling an erythroid, DNA-binding complex which includes the TAL1,
681 E47, GATA-1 and Ldb1/NLI proteins. *EMBO J*. 1997;16(11):3145–57.
- 682 59. van Galen P, Kreso A, Wienholds E, et al. Reduced Lymphoid Lineage Priming
683 Promotes Human Hematopoietic Stem Cell Expansion. *Cell Stem Cell*.
684 2014;14(1):94–106.
- 685 60. Yoshimoto M, Porayette P, Glosson NL, et al. Autonomous murine T-cell progenitor
686 production in the extra-embryonic yolk sac before HSC emergence. *Blood*.
687 2012;119(24):5706–14.
- 688 61. Gentek R, Ghigo C, Hoeffel G, et al. Epidermal $\gamma\delta$ T cells originate from yolk sac
689 hematopoiesis and clonally self-renew in the adult. *J. Exp. Med*.
690 2018;215(12):2994–3005.
- 691 62. Rybtsov S, Batsivari A, Bilotkach K, et al. Tracing the origin of the HSC hierarchy
692 reveals an SCF-dependent, IL-3-independent CD43- embryonic precursor. *Stem
693 Cell Reports*. 2014;3(3):489–501.
- 694 63. Wang C, Wang L, Su B, et al. Serine protease inhibitor Kazal type 1 promotes
695 epithelial-mesenchymal transition through EGFR signaling pathway in prostate
696 cancer. *Prostate*. 2014;74(7):689–701.
- 697 64. Chen F, Long Q, Fu D, et al. Targeting SPINK1 in the damaged tumour
698 microenvironment alleviates therapeutic resistance. *Nat. Commun*. 2018;9(1):1–19.

699

700 **Figure Legends**

701

702 **Figure 1 - Single cell transcriptome analysis reveals clusters of naïve and lineage-**
703 **committed haematopoietic progenitors.**

704 **(A)** Schematic of the single cell RNA sequencing experiment where iPSCs (SFCi55) were
705 differentiated in vitro (IVD) for 13 days (Supplementary Figure1A), CD235a⁺CD43⁺
706 suspension cells were isolated by flow cytometry and subjected to 10x genomics
707 sequencing platform. **(B)** tSNE visualization of 11,420 cells divided into 8 clusters including
708 clusters defined by gene expression as naïve (naïve 1 and 2), and others that expressed
709 genes associated with erythroid (Ery1 and 2), megakaryocyte (Mega 1,2 and 3) and
710 granulocyte (granulo) lineages. **(C)** Heatmap showing expression of the top 10 marker
711 genes for each cluster (colors for each cluster as in Figure 1B). **(D)** Gene expression
712 levels of marker genes associated with different progenitor cell types that were identified
713 by clustering, visualized on tSNE.

714

715 **Figure 2 - Trajectory analyses support naïve progenitor identity and their**
716 **progression to lineage committed progenitors.**

717 **(A)** Diffusion plot displays the naïve progenitors in the core region of the plot from where
718 the three direction of commitment originates, the arrows indicates the commitment
719 directions. **(B)** Representation of each cluster on the diffusion plot. **(C)** PAGA analysis
720 show that naïve cluster are connected to the lineage committed cells. Each node contains

721 a pie chart showing the proportion of cells for each cluster. Colors indicate cluster
 722 identities. **(D)** Monocle trajectory analyses demonstrates a similar pattern to that obtained
 723 from the diffusion plot shown in A with naïve progenitors at the top of the hierarchy, with
 724 progression toward committed. **(E)** Monocle trajectory visualizing each cluster
 725 individually. **(F)** Expression levels of the marker genes coding for transcription factors
 726 associated with each cluster, bars color indicates the cluster.
 727

728 **Figure 3 - CD44 identifies clonogenic hematopoietic progenitors.**

729 **(A)** Expression levels of genes encoding cell surface markers, *CD33*, *CD44* and *CD18* that
 730 were associated with the naïve progenitor clusters. **(B)** Scatter plot of flow cytometry
 731 profile of naïve 1A, 1B and 2 cells at day 13 of differentiation (hiPSCs-SFCi55). Cells are
 732 gated on CD235a⁻CD43⁺CD33⁺. **(C)** Schematic of the chimeric culture system using the
 733 ZsGreen reporter to trace cells during the differentiation process. ZsGreen and parental
 734 line (SFCi55) were differentiated in a synchronous manner, at day 10 naïve 1 and naïve 2
 735 cells are sorted and co-cultured with the parental line differentiation. Co-culture is then
 736 analyzed at day 13. **(D)** Representative flow cytometry profile of the day 13 naïve
 737 descendants' cells after sorting at day 10 and chimeric co-culturing of naïve 1 (teal) and
 738 naïve 2 (pink) cells. Contribution of naïve 1, in teal, and naïve 2, in pink, to the naïve 1A,
 739 1B and 2 compartment (n=6, multinomial logistic regression, *p<0.05, **p<0.01,
 740 ***p<0.005). **(E)** CFU-C analyses of FAC-sorted naïve 1 and naïve 2 cells from day 10
 741 (n=3, paired t-Test p=0.0753) (hiPSCs-SFCi55). **(F)** CFU-C analyses of FAC-sorted naïve
 742 1A, 1B and 2 cells from day 13 (n=9, Holm-Sidak's test, p<=0.001) (hiPSCs-SFCi55). **(G)**
 743 RUNX1-GFP expression in naïve 1 and naïve 2 at both day 10 and day 13 (n=12, paired t
 744 test; * p<0.05, ** p<0.01).
 745

746 **Figure 4 - CITE-seq analyses confirm markers for naïve and lineage-committed
 747 progenitor cells.**

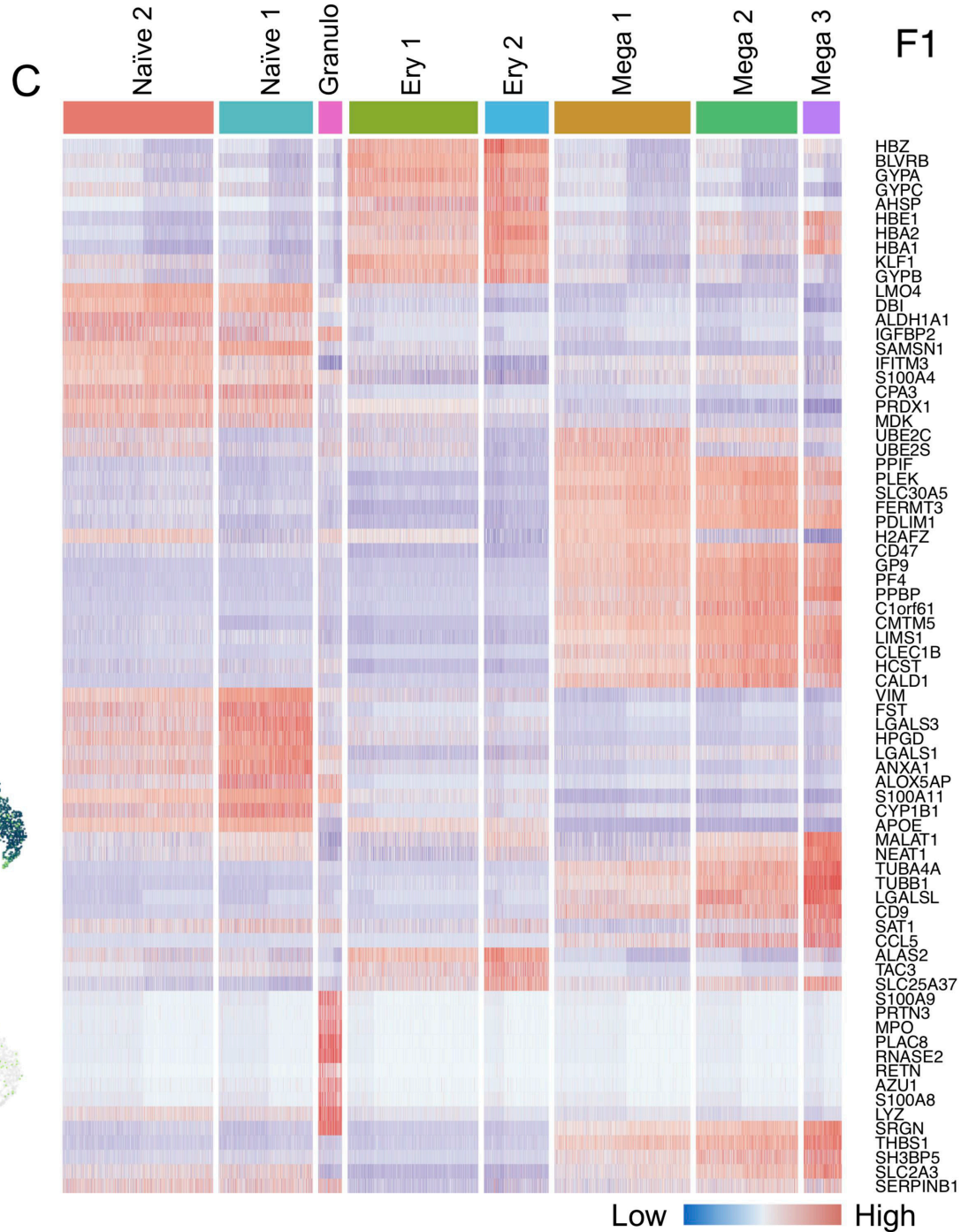
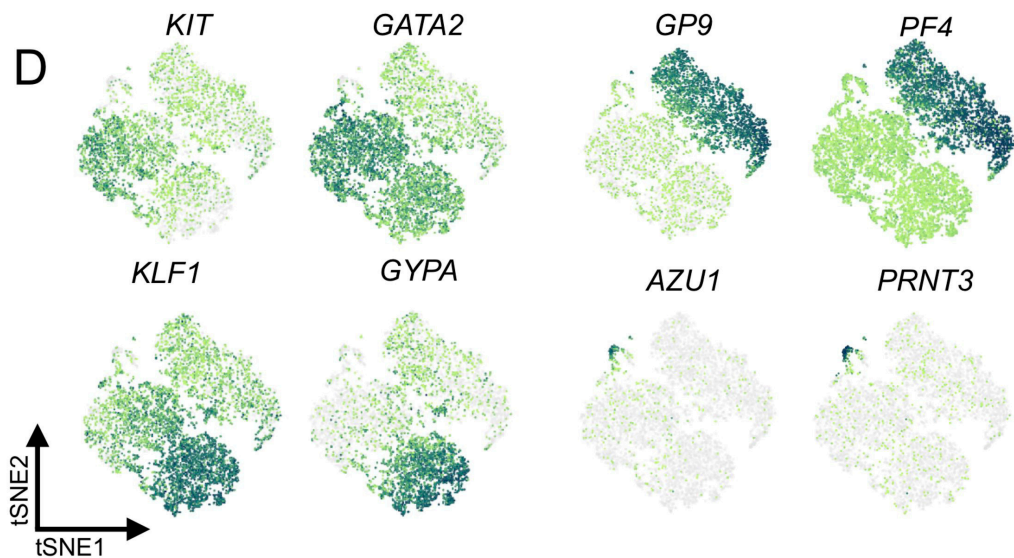
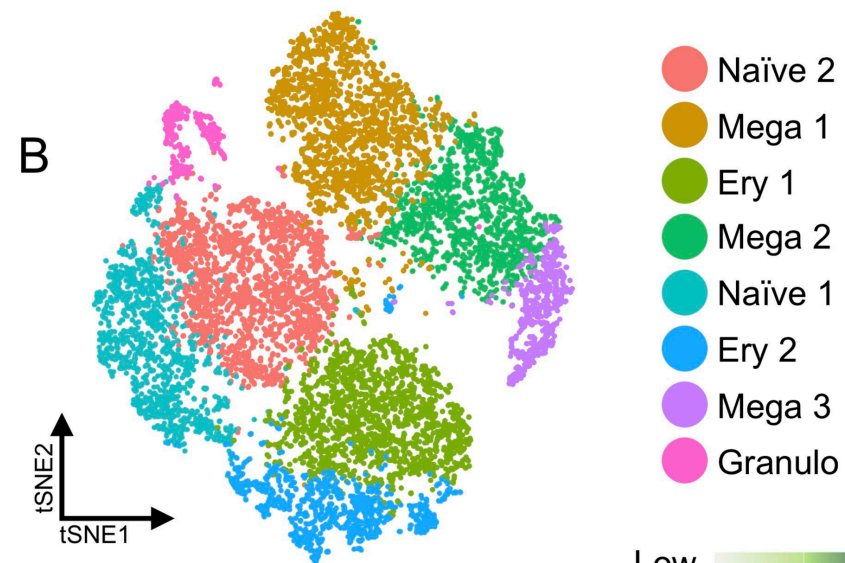
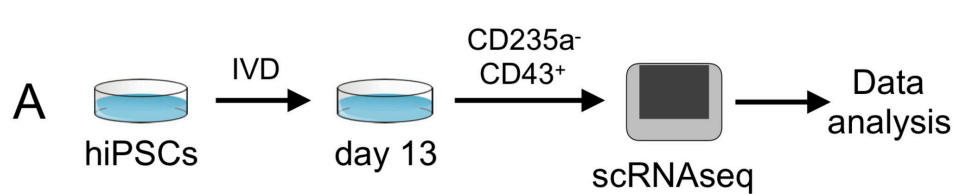
748 **(A)** tSNE visualization of the CITE-seq analysis of day 10 and day 13 CD235a⁻ suspension
 749 cells; reduction and clustering were performed using only transcriptomic data (hiPSC-
 750 SFCi55). **(B)** tSNE visualization of the libraries obtained from CD235a⁻ suspension cells
 751 collected at day 10 (pink) and day 13 (teal) showing lineage commitment direction. **(C)**
 752 Single cell protein expression level of the membrane markers associated with the different
 753 cell types. Data are visualized on tSNE (ADT = antibody derived tags). **(D)** tSNE plot and
 754 annotation of clustering obtained from analysis derived from ADT data alone. **(E)**
 755 Visualization of the libraries obtained from CD235a⁻ suspension cells projected on the
 756 tSNE obtained from ADT data in D, cell progression shows lineage commitment trajectory.
 757 Cells are colored according to the day of collection (day 10 in pink and day 13 in teal). **(F)**
 758 Confusion matrix of clustering obtained from complete transcriptomic data (RNA) and that
 759 obtained from ADT data alone. Color of each box indicates the % of cells classified in each
 760 RNA versus ADT cluster. **(G)** Gene expression levels of human foetal gene marker genes
 761 in our naïve and lineage committed progenitors.
 762

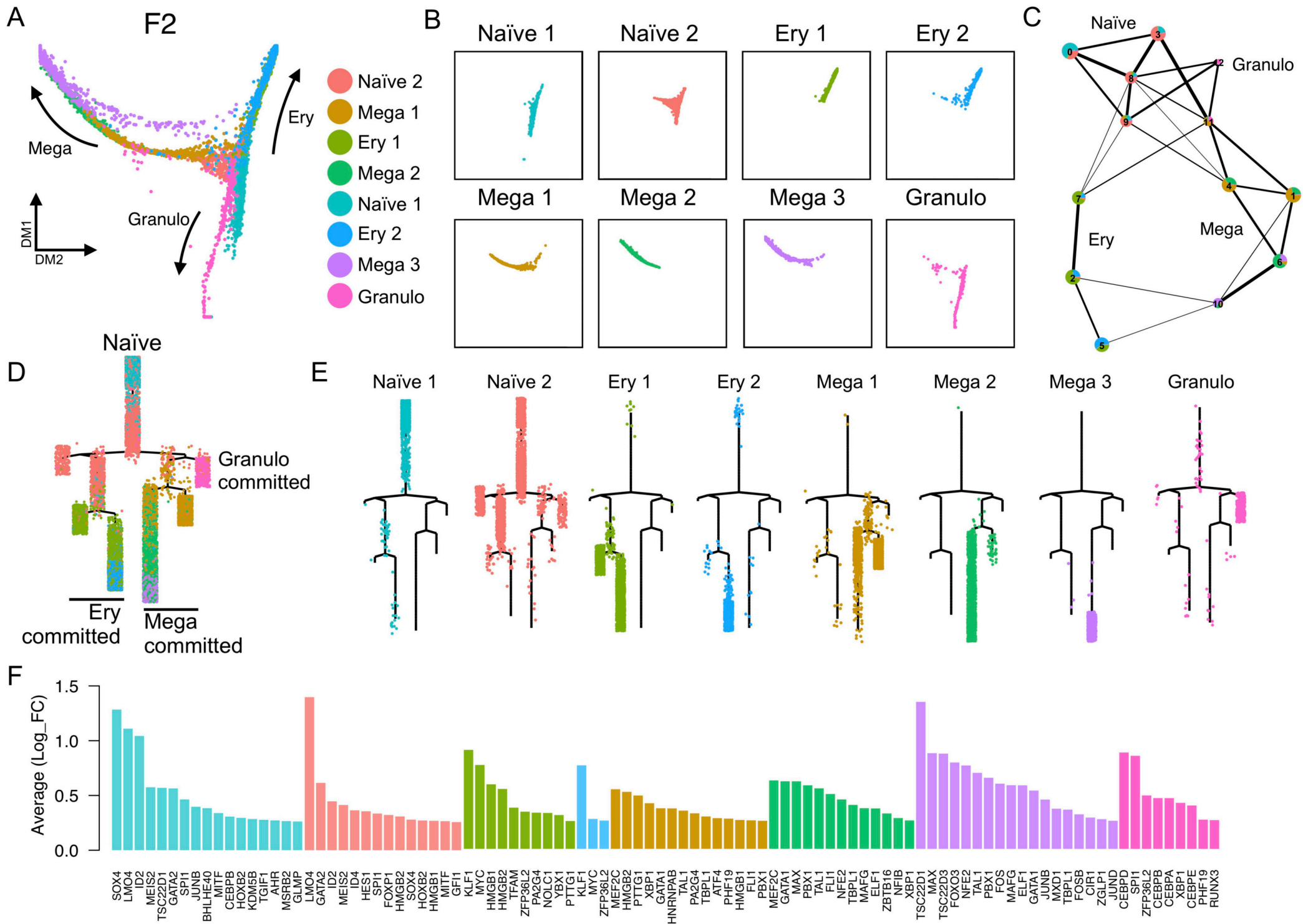
763 **Figure 5 - Artificial neural network identifies HSC-like cells in iPSC derived**
764 **hematopoietic cells.**

765 **(A)** Schematic of the artificial neural network (ANN) architecture for label-transfer. An
766 input layer (3479 units), two fully connected hidden layers (64 and 32 units) and a 28-unit
767 softmax layer corresponding to cell types in the source domain (human foetal liver
768 scRNAseq data) used for training. Classification of cell types in the target domain of
769 human iPSC-derived single cell transcriptomes (test data). **(B)** Proportion of cells labelled
770 HSC/MPPs at day 10 or day 13 of hiPSC differentiation in vitro. **(C)** Proportion of in vitro
771 derived CD235a⁻ progenitors and CD235a⁻CD43⁺ cells labelled HSC/MPP by the ANN (ND
772 = not detected). **(D)** UMAP visualization of the integrated dataset containing in vivo derived
773 (blue) and in vitro annotated (pink) HSC/MPPs. **(E)** Heatmap of differentially expressed
774 genes coding for transcription factors obtained from the comparison of in vivo and in vitro
775 derived HSC/MPPs. **(F)** Transcription factor binding motifs enriched upstream of the 54
776 genes identified as differentially expressed between HSC/MPP generated in vitro and in
777 vivo.

778

779





F3

