

Single-cell RNA-Sequencing uncovers transcriptional states and fate decisions in haematopoiesis

Emmanouil I. Athanasiadis¹⁻³, Helena Andres^{4,*}, Jan G. Botthof^{1-3,*}, Lauren Ferreira¹⁻³, Pietro Lio⁴, Ana Cvejic¹⁻³

¹Department of Haematology, University of Cambridge, Cambridge, CB2 0XY, UK

²Wellcome Trust Sanger Institute, Wellcome Trust Genome Campus, Cambridge, CB10 1SA, UK

³Wellcome Trust – Medical Research Council Cambridge Stem Cell Institute, Cambridge, CB2 1QR, UK

⁴Computer Laboratory, University of Cambridge, Cambridge, CB3 0FD, UK

Correspondence: Ana Cvejic, Wellcome Trust Sanger Institute, Wellcome Genome Campus, The Morgan Building, Hinxton, Cambridge, Cambridgeshire, CB10 1SA.

E-mail: as889@cam.ac.uk

*H.A. and J.G.B contributed equally to this study

ABSTRACT

The success of marker-based approaches for dissecting haematopoiesis in mouse and human is reliant on the presence of well-defined cell-surface markers specific for diverse progenitor populations. An inherent problem with this approach is that the presence of specific cell surface markers does not directly reflect the transcriptional state of a cell. Here we used a marker-free approach to computationally reconstruct the blood lineage tree in zebrafish and order cells along their differentiation trajectory, based on their global transcriptional differences. Within the population of transcriptionally similar stem and progenitor cells our analysis revealed considerable cell-to-cell differences in their probability to transition to another, committed state. Once fate decision was executed, the suppression of transcription of ribosomal genes and up-regulation of lineage specific factors coordinately controlled lineage differentiation. Evolutionary analysis further demonstrated that this haematopoietic program was highly conserved between zebrafish and higher vertebrates.

INTRODUCTION

Mammalian blood formation is the most intensely studied system of stem cell biology, with the ultimate aim to obtain a comprehensive understanding of the molecular mechanisms controlling fate-determining events. A single cell type, the haematopoietic stem cell (HSC), is responsible for generating more than 10 different blood cell types throughout the lifetime of an organism¹. This diversity in the lineage output of HSCs is traditionally presented as a step wise progression of distinct, transcriptionally homogeneous populations of cells along a hierarchical differentiation tree²⁻⁶. However, most of the data used to explain the molecular basis of lineage differentiation and commitment were derived from populations of cells isolated based on well-defined cell surface markers⁷. One drawback of this approach is that a limited number of markers is used simultaneously to define the blood cell identity. Consequently, only a subpopulation of the overall cellular pool is examined and isolated cells, although homogeneous for the selected markers, show considerable transcriptional and functional heterogeneity⁸⁻¹². This led to the development of various refined sorting strategies in which new combinations of marker genes were considered to better “match” the transcriptional and functional properties of the cells of interest.

The traditional model of haematopoiesis assumes a stepwise set of binary choices with early and irreversible segregation of lymphoid and myeloid differentiation pathways^{2, 3}. However, the identification of lymphoid-primed multipotent progenitors (LMPP)⁴, which have granulocytic, monocytic and lymphoid potential but low potential to form megakaryocyte and erythroid lineages prompted development of alternative models of haematopoiesis. More recently, it has been demonstrated that megakaryocyte-erythroid progenitors can progress directly from HSC without going through a common myeloid intermediate (CMP)¹³; or that the stem cell compartment is multipotent, while the progenitors are unipotent⁶. Clear consensus on the lineage branching map, however, is still lacking.

Recent advances in single-cell transcriptional methods have made it possible to investigate cellular states and their transitions during differentiation, allowing elucidation of cell fate decision mechanisms in greater detail. Computational ordering methods have proved to be particularly useful in reconstructing the differentiation process based on the transcriptional changes of cells at different stages of lineage progression¹⁴⁻¹⁶.

Here we created a comprehensive atlas of single cell gene expression in adult zebrafish blood cells and computationally reconstructed the blood lineage tree *in vivo*. Conceptually, our approach differs from the marker based method in that the identity of the cell type/state is determined in an unbiased way i.e. without prior knowledge of surface markers. The

transcriptome of each cell was projected on the reconstructed differentiation path giving complete insight into the cell state transitions occurring during blood differentiation. Importantly, development of this strategy allowed us, for the first time, to assess haematopoiesis in a vertebrate species in which surface marker genes/antibodies are not readily available. Finally, this study provides unique insight into the regulation of haematopoiesis in zebrafish and also, along with complementary data from mouse and human, addresses the question of interspecies similarities of haematopoiesis in vertebrates.

RESULTS

Single cell RNA-Sequencing analysis of 1,422 zebrafish haematopoietic cells

As an alternative to marker-based cellular dissection of haematopoietic hierarchy, we have set out to classify haematopoietic cells based on their unique transcriptional state. We started by combining FACS index sorting with single cell RNA-Seq to reveal the cellular properties and gene expression of a large number of blood cells simultaneously. To cover the entire differentiation continuum, kidney derived blood cells from eight different zebrafish transgenic reporter lines and one non-transgenic line were FACS sorted (Fig. 1a, Supplementary Table 1). Each blood cell was collected in a single well of a 96-well plate. At the same time, information about the cell size (FSC) and granularity (SSC), as well as the level of the fluorescence, were recorded.

RNA from each cell was isolated and used to construct a single mRNA-Seq library per cell, which was then sequenced to a depth of around 1×10^6 reads per library. Following quality control (QC) 1,422 cells were used for further analysis (Supplementary Fig. 1 and 2a-d). Importantly, the average single-cell profiles showed good correlation with independent bulk samples ($PCC=0.7-0.9$, Supplementary Fig. 2e). In addition, PCA, ICA and Diffusion maps (Supplementary Fig. 3a) showed that cells were intermixed irrespective of the fish or the plate they originated from. This confirmed that the cells were separated in the analyses based on their biological differences rather than batch induced biases.

HSPC can reach specific cell fates through a single path in the “state-space”

A dynamic repertoire of gene expression in thousands of cells during differentiation could be used to infer a single branched differentiation trajectory. Due to the unsynchronised nature of haematopoiesis each single cell exhibits a different degree of differentiation along the differentiation continuum. Therefore, the generated trajectory could be used to infer the differentiation path of a single cell. To examine the transcriptional transition undergone by differentiating cells, we identified the 1,845 most highly variable genes (Supplementary Fig.

3b) and performed expression based ordering using Monocle2¹⁵. Based on global gene expression profiles of the cells, we identified five (1-5) distinct cell “states” (Fig. 1b).

Differential expression analysis of each state versus all other states, followed by gene ontology (GO) enrichment analysis (see methods), provided clear insights into the cell types in each state (Fig. 1c). Specifically, state 1 contains GO terms relating to antigen processing, including genes that are highly expressed in the monocyte lineage, such as *cd74a/b*¹⁷, *ctss2.2*¹⁸ and *mhc2dab*¹⁹ (Supplementary Table 2). The functionality of state 2 relates to leukocyte migration, including genes specific to neutrophils (e.g. *cxcr4b*²⁰, *rac2*²¹ and *wasb*^{22, 23} (Supplementary Table 2). State 3 is highly enriched for genes that are involved in ribosome biogenesis, including *fbf* (Fibrillarin) and *pes* (Pescadilo), both of which are critical for stem cell survival^{24, 25} (Supplementary Table 2). Since there is also enrichment for HSC homeostasis, this state is most likely to be haematopoietic stem/progenitor cells (HSPCs). With GO terms that include gas exchange and erythrocyte differentiation involving the adult haemoglobins, *ba1*, *ba1l* and *hbaa1*²⁶ together with the erythroid-specific aquaporin gene, *aqp1a*²⁷ (Supplementary Table 2), state 4 can be assigned to the erythroid lineage. Finally, state 5 has functionality that is relevant for circulatory system development and blood coagulation, both of which include *itga2b* (also known as *cd41*) together with its heterodimer *itgb3b*²⁸ (Supplementary Table 2). Since these gene lists include other genes that interact with this platelet integrin receptor complex, as well as additional genes relevant for platelet function, we assigned this cell state to thrombocytes.

To experimentally confirm our computational predictions, we sorted cells from transgenic lines that were the most abundant in each of the five states (Fig. 2) and stained them using May-Grünwald Giemsa staining. Indeed, the morphological properties of the sorted cells (Fig. 1c, Supplementary Fig. 4-5) matched the assigned cell types, therefore adding confidence to these cell type assignments. As expected, the signature genes such as *marco*, *lyzC*, *hhex*, *alas2* and *itga2b* were within the most differentially expressed genes in monocytes, neutrophils, HSPC, erythrocytes and thrombocytes respectively (Fig. 1d).

Taken together, the reconstructed branched tree revealed a gradual transition of myeloid cells from immature to more differentiated cells. Within this tree, HSPCs assumed a new committed state through a single path, suggesting that during steady state haematopoiesis, HSPCs can reach a specific cell fate through only one type of intermediate progenitor.

Cells within distinct states differ in their repopulation potential

Functional *in vivo* transplantation assays have been traditionally used to assess the differentiation potential of different haematopoietic populations. To examine the repopulation and lineage potential of the cells within different states we sorted cells from *Tg(mpx:EGFP)*²⁹, *Tg(gata1:EGFP)*³⁰ and *Tg(runx1:mCherry)*³¹ fish to enrich for neutrophil, erythroid and HSPC cell state respectively. We next injected 500 donor cells into sub-lethally irradiated, immunocompromised *rag2*^{E450fs/-} zebrafish³² and assessed their engraftment at one day, four- and fourteen weeks post injection (PI) (Fig. 3a).

Analysis of kidney repopulation revealed that *mpx+*, *gata1+* and *runx1+* cells were able to home to the kidney oneday PI (Fig. 3b). However, only progeny of *runx1+* cells were detectable at four weeks PI in all examined recipients (Fig. 3b). No progeny of *mpx+* and *gata1+* were evident at the same time point. To examine the lineage output of *runx1+* cells following transplantation we sorted engrafted *runx1+* kidney cells four weeks PI and processed them for scRNA-Seq analysis. The scRNA-Seq data from 149 engrafted *runx1+* cells projected onto a Monocle trajectory revealed the multilineage potential of donor *runx1* cells (Fig. 3c). Importantly, the donor-derived *runx1+* cells were capable of long-term (three months) generation of blood cells in irradiated hosts, suggesting that at least some of these cells were HSCs.

According to transplantation assays, cytopins and transcriptional profiling of cells prior and following transplantation, cells located in the branches of the Monocle tree show progression of lineage restricted progenitors to mature blood cells with no repopulation potential. However, cells in the middle of the Monocle tree (state 3) are a mixture of progenitors and HSCs with long term multilineage potential.

Transcriptional changes at the branching point

To examine the heterogeneity within the HSPC population in more detail, we used a Hidden Markov model whose parameters were estimated by a Deep Neural Network (DNN) (see Methods). DNN was trained using four data sets, namely, cells from erythrocytes, thrombocytes, monocytes and neutrophil branches. Since relative weightings of the individual genes were determined from the neural network, the most “influential” genes for each of the branches were identified in an unbiased way. This allowed us to determine the probability of each individual HSPC to transition to any of the four given branches (Fig. 4) without *a priori* knowledge of which genes should be included in the analysis. Interestingly, cells that were transcriptionally similar overall and were close together on the pseudotime trajectory displayed different probabilities of transitioning to another state (Fig. 4). This

suggested that although global transcriptional changes before and after branching point were continuous, the cell fate decision itself was influenced by a subset of highly relevant genes and could not be implied based on the position of the cell on the pseudotime trajectory.

Suppression of transcription of ribosomal genes and up-regulation of lineage specific factors co-ordinately control lineage differentiation

Differentiation generally involves specific regulated changes in gene expression. To understand the dynamics of transcriptional changes during the differentiation of myeloid cells, we examined trends in gene expression in each of the four branches (Fig. 5). Dynamically expressed genes within each of the branches showed two main trends (see methods). These included genes gradually upregulated through pseudotime and genes gradually down-regulated (Fig. 5a-b).

Genes upregulated in pseudotime included well known genes related to the specific function of the relevant cell type (Fig. 5b). The majority of cells characterised as erythroid dynamically expressed genes such as *alas2*, *aqp1a.1*, *ba1*, *ba1l*, *cahz* and *hbaa1*. Similarly, cells in the monocyte branch dynamically expressed genes like *c1qa*, *cd74a*, *ifngr1*, *marco*, *myod1* and *spi1a*; among other genes the *cebpb*, *cfl1*, *cxcr4b*, *illr4*, *mpx* and *ncf1* were upregulated in pseudotime in the neutrophil branch and thrombocytes dynamically expressed *fn1b*, *gp1bb*, *itga2b*, *mpl*, *pbx1a* and *thbs1b*. A complete list of all genes that are dynamically expressed across pseudotime can be found in Supplementary Table 2.

Interestingly, genes downregulated through pseudotime (Fig. 5b) in each of the four branches were consistently enriched for genes involved in ribosome biosynthesis, as revealed by GO terms “biosynthetic process”, “ribosome” and “translation” (Supplementary Table 2). This is an interesting finding, because previous studies suggested that HSCs have significantly lower rates of protein synthesis than other haematopoietic cells³³. Therefore, we went on to investigate the expression of ribosomal proteins in pseudotime in greater depth (Fig. 5c).

Out of 168 genes annotated as “ribosomal proteins” on Ensembl BioMart database (Supplementary Table 2), 89 genes had low, random expression in our dataset (Fig. 5c). These genes encoded mainly mitochondrial ribosomal proteins (Fig. 5c). In contrast, 79 genes that showed high expression across all cells encoded cytoplasmic ribosomal proteins and were downregulated in pseudotime in all four branches (Fig 5c). These findings further indicate that there is a common developmental event in which suppression of transcription of

ribosomal genes and up-regulation of lineage specific factors direct lineage commitment and terminal differentiation.

Zebrafish have a highly conserved HSPC transcriptome compared to mouse and human

Zebrafish is an important model system in biomedical research and has been extensively used for the study of haematopoiesis. Although it has been demonstrated that many transcription factors and signaling molecules in haematopoiesis are well conserved between zebrafish and mammals³⁴, comparative analysis of the whole transcriptome was lacking.

In order to explore the evolution of blood cell type specific genes, we performed conservation analysis between zebrafish and other vertebrate species (see Methods). For this analysis, we enriched our initial dataset with 81 natural killer (NK) and 109 T-cells derived from the spleen of two adult zebrafish³⁵. Our analysis revealed particularly high conservation of the HSPC transcriptome. For example, 90% of HSPC specific genes in zebrafish had an ortholog in human and mouse compared to 70-80% of erythrocyte-, monocyte-, neutrophil- and thrombocyte-specific genes (Fig. 6a). The lowest conservation was observed for T-cells (59%) and NK cells (68%), possibly reflecting their adaptation to fish specific pathogens and virulence factors (Fig. 6a).

Gene duplication is the major process of gene divergence during the molecular evolution of species³⁶. We therefore analysed duplications that occurred exclusively before (referenced hereafter as pre-speciation genes) or after speciation (referenced hereafter as post-speciation genes) of the last common ancestor between fish (Actinopterygii) and mammals (Sarcopterygii)^{35, 37}, (see methods section). Out of 7,424 paralogs that were expressed in our data set (see Methods) around 79% were duplicated pre- and 21% were duplicated post-speciation (Fig. 6b). Following ray-finned specific duplication, the paralogs were more likely to functionally diverge (88%) and show expression in different cell types than to remain expressed in the same cell type (conserved expression), 12% (Fig. 6b and c). Interestingly, HSPCs had the highest percentage of paralogs (19%) with a conserved expression pattern (Fig. 6c). This number was lowest for duplicated genes in innate (0% for the neutrophils and 6% in monocytes) and adaptive immune cells (8% for the NK and 6% for the T-cells). Altogether our findings further underline the relevance of the zebrafish model system in advancing our understanding of the genetic regulation of haematopoiesis in both normal and pathological states.

BASiCz - Blood Atlas of Single Cells in zebrafish

The characterisation of mouse and human haematopoietic cells is dependent on the presence of cell-surface markers and availability of antibodies specific for diverse progenitor populations. The antibodies for these cell surface markers are thus used to isolate relatively homogeneous cell populations by flow cytometry. Transcriptional profiling of isolated cell populations³⁸⁻⁴⁰ and more recently single cells⁴¹, have further allowed genome-wide identification of cell-type specific genes. However, beyond mouse and human, less is known about the transcriptome of blood cell types, mainly due to the lack of suitable antibodies.

To overcome this knowledge gap, we have generated a user-friendly cloud repository, BASiCz (Blood Atlas of Single Cells in zebrafish) for interactive exploration and visualisation of 31,953 zebrafish genes in 1,422 haematopoietic cells across five different cell types. The generated database (<http://www.sanger.ac.uk/science/tools/basicz>) allows easy access and retrieval of sequencing data from zebrafish myeloid cells.

DISCUSSION

Cell differentiation during normal blood formation is considered to be an irreversible process with a clear directionality of progression from HSCs to more than 10 different blood cell types. It is, however, widely debated to what extent the process is gradual or direct^{6, 13} on the cellular level; and in the case of the gradual model, what the intermediates of the increasingly restricted differentiation output of progenitor cells are²⁻⁵. Although these models are very different in the way that they describe lineage progression, the identity of haematopoietic cells is determined based on the cell surface markers and the progression of cells during differentiation is defined on a cellular rather than transcriptional level.

Here we used a marker free approach to order cells along their differentiation trajectory based on the transcriptional changes detected in the single cell RNA-Seq dataset. Our analysis showed a gradual transition of cells on a global transcriptional level from multipotent to lineage restricted. The computationally reconstructed tree further revealed that differentiating cells moved along a single path in the “state-space”. This path included an early split of cells towards thrombocyte-erythrocyte and monocyte-neutrophil trajectories. However, cells in the “middle” of the tree (HSPC state) showed considerable cell-to-cell variability in their probability to transition to any of the four cell types. This suggested that although global transcriptional changes before and after the branching point were continuous, the probability of a cell transitioning to any of the four committed states was determined only by a subset of highly relevant genes. Therefore, cells that were

transcriptionally similar overall could have a high probability of differentiation to distinct cell types.

Interestingly, once the cell fate decision was executed, suppression of transcription of ribosomal genes and up-regulation of genes which are relevant for the function of each cell type coordinately controlled lineage differentiation. Of all genes that were annotated as “ribosomal proteins” on the Ensembl BioMart database, only those that encoded cytoplasmic ribosomal proteins showed dynamic expression in pseudotime in our dataset. These findings are not in line with previous studies, which suggested that HSCs have significantly lower rates of protein synthesis compared to other haematopoietic cells. It should be noted, however, that in this study we measured the transcription of genes that encoded ribosomal proteins rather than *de novo* protein synthesis like in³³. Thus, one plausible explanation for the observed discrepancies is a low correlation between transcription of the ribosomal genes and protein production and that these two processes are to some extent uncoupled during blood differentiation.

Our comparative analysis between zebrafish, mouse and human across seven different haematopoietic cell types revealed a high overall conservation of blood cell type specific genes. Together with BASiCz, a user-friendly cloud repository, we generated a comprehensive atlas of single-cell gene expression in adult zebrafish blood. Data-driven classification of cell types provided high-resolution transcriptional maps of cellular states during differentiation. This allowed us to define the haematopoietic lineage branching map, for the first time, in zebrafish *in vivo*.

METHODS

Zebrafish Strains and Maintenance

The maintenance of wild-type (Tubingen Long Fin) and transgenic zebrafish lines^{29-31, 42-46} (Supplementary Table 1) was performed in accordance with EU regulations on laboratory animals, as previously described⁴⁷.

Single-Cell Sorting

A single kidney from heterozygote transgenic or wild-type fish was dissected and placed in ice cold PBS/5% fetal bovine serum. At the same time, testes were dissected from the same fish. Single cell suspensions were generated by first passing through a 40 um strainer using the plunger of a 1 ml syringe as a pestle. These were then passed through a 20 um strainer before adding 4',6-diamidino-2-phenylindole (DAPI, Beckman Coulter, cat no B30437) for *mCherry/dsRed2*, or propidium iodide (PI, Sigma cat no P4864) for *GFP/EGFP*. Individual

cells were index sorted into wells of a 96 well plate using a BD Influx Index Sorter. Kidneys from a non-transgenic line were used as a control for gating¹⁶.

Whole Transcriptome Amplification

The Smart-seq2 protocol^{48, 49} was used for whole transcriptome amplification and library preparation as described previously¹⁶ using 92 External RNA Controls Consortium (ERCC) spike-ins at a final dilution of 1:10⁷. These were sequenced on the Illumina Hi-Seq2500 or Hi-Seq4000 platforms.

Cytology

Sorted transgene-positive or gated wild type cells were concentrated by cytocentrifugation at 350 rpm for 5 minutes onto SuperFrostPlus slides using a Shandon Cytospin 3 cytocentrifuge. Slides were fixed for 3 minutes in -20 °C methanol and stained with May-Grünwald Giemsa (Sigma) as described elsewhere⁵⁰. Images were captured as described elsewhere⁴⁷.

Transplantation experiments

Adult *rag2*^{E450fs/-} mutant fish³² were irradiated in an IBL 437 irradiator using a 10 Gy dose from a Caesium 137 source. After 1-2 days of recovery, donor cells were prepared from kidneys of transgenic fish as described above. Using the same gating strategy as employed for the single cell sorting, fluorescent cells were collected by flow cytometry into microtubes containing 20 ul ice cold PBS/5% fetal bovine serum. Using a volume of 10 ul, 500 cells were transplanted into the anaesthetised (0.02% tricaine, Sigma A5040) *rag2*^{E450fs/-} recipients via intraperitoneal injection. As described above, engraftment into the whole kidney marrow was analysed by FACS at one day, four- and fourteen weeks post transplantation. The engrafted cells at four weeks post transplantation were single cell index sorted and processed for single cell RNA-Seq as described above.

Single cell RNAseq processing and Quality Control

Reads were aligned to the zebrafish reference genome (Ensemble BioMart version 83) combined with the *EGFP*, *mCherry*, *tdTomato* and ERCC spike-ins sequences. Quantification was performed using Sailfish⁵¹ version 0.9.0 with the default parameters using paired-end mode (parameter -l IU).

Transcript Per Million (TPM) values reported by Sailfish were used for the quality control (QC) of the samples. Wells with fewer than 1,000 expressed genes (TPM>1), or more than 60% of ERCC or Mitochondrial content were initially annotated as poor quality cells

(Supplementary Fig. 1). However, due to the lower number of expressed genes in erythroid cells, we further investigated the expression levels of adult globin genes, *ba1* and *hbaa1*²⁶, in all erythroid cells. Based on comparison with the empty wells, samples that expressed both *ba1* (> 40,000 TPM) and *hbaa1* (> 9000 TPM) were considered to pass QC (Supplementary Fig. 2). Therefore, a total of 1,422 single cells were selected for further analysis.

Average single-cell profiles compared to corresponding bulk wells revealed strong correlations (Pearson's Correlation Coefficient) ranging from 0.7 to 0.9 as illustrated in Supplementary Fig. 2, suggesting that the single cell expression profiles were effectively quantified.

For each of the 1,422 single cells, both gene and ERCC counts reported by the Sailfish, were transformed into normalised counts per million (CPM). The library size and cell-specific biases were removed (e.g. differences during amplification, ERCC concentration, batch effects etc.) using the scran R package (version 1.3.0) published in⁵². Out of 31,953 genes, we retained those that were expressed in at least 1% of all cells (CPM>1). Thus, a total of 20,960 genes were used for further analysis.

Technical noise fit and identification of highly variable genes

To distinguish biological variability from the technical noise in our single-cell experiments we inferred the most highly variable genes using ERCCs as spike-in in all 1,422 blood cells⁵³. We used the scLVM⁵⁴ R package (version 0.99.2) to identify the 1,845 most highly variable genes (Supplementary Fig. 3).

Principal Component Analysis (pcaMethods (version 1.64.0)), Independent Component Analysis (FastICA (version 1.2) and Diffusion Maps (destiny⁵⁵ (version 1.3.4)), were used to verify that all cells were intermixed in the reconstructed 3D component space based on their transcriptional properties and not based on the fish or a plate they originated from.

Pseudotime ordering of zebrafish haematopoietic cell, differential expression analysis and the analysis of dynamically expressed genes

The set of 1,845 most highly variable genes was used to order the 1,422 single cells along a trajectory using the Monocle2¹⁵ R package (version 1.99.0). The “tobit” expression family and “DDRTree” reduction method were used with the default parameters. As illustrated in Fig. 1, cells ordered in the pseudotime created five distinct states. To assign identity to each of the five states, we performed differential expression (DE) analysis between each state

versus the remaining four using the “differentialGeneTest” Monocle2 function. We modeled expression profiles of each state using a Tobit family generalized linear model (GLM) as described in¹⁵. For each state, statistically significant genes that scored $P < 0.01$, $q < 0.1$ (False Discovery Rate) and were expressed in more than 50% of the cells were further used to perform Gene Ontology (GO) analysis.

Finally, we identified genes that change as a function of pseudotime across each of the four branches by setting the “fullModelFormulaStr” parameter equal to “~sm.ns(Pseudotime)”. Genes whose expression changed dynamically in pseudotime were selected using the same statistical criteria as described for DE genes. For each branch we clustered dynamically expressed genes using the “plot_pseudotime_heatmap” function with the default parameters. The number of clusters (trends) in each branch was determined by its silhouette plot score (cluster R package version 2.0.5). To generate the trend lines across different states (see Fig. 3b), we used the average expression pattern of the dynamically expressed genes that follow the same trend across pseudotime and fit them using *ggplot2* R package (version 2.2.1) *stat_smooth()* parameter. We used the Gaussian linear model and formula the “ $y \sim poly(x,2)$ ” at 0.95 of standard error (gray area of the plot).

For the analysis of ribosomal genes, we used the Ensembl BioMart version 83 and selected all genes annotated with the term “ribosomal protein”. We performed clustering using the pheatmap function (R pheatmap package version 1.0.8) using Euclidean distance and ward.D2 linkage.

Gene Ontology (GO) analysis

DE genes were ranked for each of the five states based on the mean \log_{10} counts. Genes with average lower than 2 and those expressed in more than one state were not included in the GO analysis. GO analysis was performed using the gProfileR⁵⁶ package (Version 0.6.1) using the gprofiler command with the following parameters: organism = ‘drerio’, hier_filtering = ‘moderate’, correction_method = ‘fdr’ and max_p_value = 0.05.

Conservation analysis of the cell type specific genes in zebrafish

In order to perform the conservation analysis, we identified the orthologous genes (BioMart Ensembl Version 83) between the zebrafish and other vertebrate species, including cave fish, tilapia, amazon molly, tetraodon, fugu, cod, human, chimpanzee, mouse, rat, dolphin, wallaby, chicken, lizard, *Xenopus*, coelacanth and lamprey. For this analysis, we enriched our initial dataset with 81 natural killer (NK) and 109 T-cells derived from the spleen of two adult zebrafish³⁵. Following the same computational approach as we did with the initial

dataset, we re-calculated the DE genes for each of the seven different clusters. We only considered “protein_coding” genes that were expressed in more than 50% of cells within each cluster and scored more than mean \log_{10} counts. This resulted in 41 erythrocyte-, 113 monocyte-, 102 neutrophil-, 212 thrombocyte-, 60 HSPC-, 34 NK- and 34 T- specific genes that were used for the further analysis. For the case of the non-DE genes, we included only “protein_coding” annotated genes that were expressed in more than 1% of all cells (CPM>1) and with average gene expression higher than the global mean of 0.10. The final list of the non-DE genes included 8,127 genes.

Analysis of duplicated genes in zebrafish

In order to analyse duplicated genes³⁵, we first identified all zebrafish “protein_coding” paralog genes listed in Ensembl (BioMart Ensembl Version 83) and split them into two groups: 1) 17,158 pre ray-finned fish duplicated genes, including *Euteleostomi*, *Bilateria*, *Chordata*, *Vertebrata* and *Opisthokonta* parent taxa, and 2) 11,806 post ray-finned fish duplicated genes, including *Neopterygii*, *Otophysa*, *Clupeocephala* and *Danio rerio* children taxa. We next removed duplicated genes that were found in common between the two groups. This resulted in 8,601 pre-, and 3,249 post-ray-finned fish genes that we used in further analysis.

For the analysis of the expression pattern divergence, we focused on genes that were expressed in our data set. We analyzed expression pattern of all paralogs of DE genes (i.e. erythrocytes, monocytes, neutrophils, thrombocytes, HSPCs, NK- and T cells) that were expressed in more than 10% of cell in each of the branches (cell states). The expression pattern was considered to be conserved if duplicated genes and their annotated paralogs were all expressed in the same cell type. However, if at least one of the paralogs was expressed in a different cell type, this was considered as an example of potential functional divergence.

Deep Neural Network (DNN) Classifier

To generate the DNN model we used Keras, a python based Deep Learning Library for Theano and Tensorflow. We worked with the Keras functional API, which allows the definition of complex systems such as multi-output models.

The DNN was used to predict the probabilities of a specific Gene Expression profile to be classified into one of the four differentiated cell types. We used the entire set of genes for all differentiated cells in the branches (1,177 cells in total) i.e. erythrocytes, thrombocytes, neutrophils and monocytes. The input is therefore formed by 31,953 nodes (genes) which

are normalized using z-values or standard score. For the hyper-parametric fine tuning of the DNN we generated and evaluated models with different number of hidden layers, hidden nodes, network initializations, regularizations and batch normalization. The final hyper parameters were chosen according to the optimal performance and convergence of the accuracy and loss values. The model was comprised of 2 hidden layers with 1000 and 500 nodes, using a weight decay regularization with a λ -value of 0.001, 'softmax' activation and 'Adam optimiser' as our optimization algorithm. The validation was performed over 10% of the initial dataset, using 'categorical cross-entropy' loss. The average classification accuracy after convergence was 0.9924 ± 0.0001 , and cross entropy loss of 0.0429 ± 0.0003 , validation accuracy of 0.9661 ± 0.0002 and cross entropy validation loss 0.1134 ± 0.0015 .

Cloud Repository

We have generated a cloud repository to enable research community to access single cell gene expression profiles of 1,422 zebrafish blood cells across all the 31,953 zebrafish genes. The implementation of the cloud service was performed using shiny (version 0.14.2) <https://shiny.rstudio.com>, and plotly (version 4.5.6) <https://plot.ly> R packages.

Statistics and reproducibility of experiments

Statistical tests were carried out using R software packages as indicated in the figure legends and the Methods section. No statistical method was used to predetermine sample sizes. Pearson Correlation Coefficient was used to compare the average profiles of single cells against the bulk. Significance of Differentially Expressed genes was calculated with an approximate likelihood ratio test (Monocle2 differentialGeneTest() function) of the full model "~state" cells against the reduced model "~1". For the Dynamically expressed genes, the full model "~sm.ns(Pseudotime)" was tested against the reduced model of no pseudotime dependence. In both cases, *P* values were normalised using the the Benjamini-Hochberg FDR (False Discovery Rate), selecting statistically significant genes with $P < 0.01$ and FDR < 0.1 . For the GO analysis, the Hypergeometric Test (equivalent to the one tailed Fisher's exact test) was used to evaluate the significant terms, while *P* values were corrected for multiple testing using the FDR approach, with FDR < 0.05 considered statistically significant, using gprofiler R package.

DATA AVAILABILITY

Raw data can be found under the accession number E-MTAB-5530 on ArrayExpress (<https://www.ebi.ac.uk/arrayexpress/>). Additional RNAseq data that were used in the present study can be found in E-MTAB-4617 and E-MTAB-3947.

ACKNOWLEDGEMENTS

The study was supported by Cancer Research UK grant number C45041/A14953 (to A.C. and E.A.), European Research Council project 677501 – ZF_Blood (to A.C.) and a core support grant from the Wellcome Trust and MRC to the Wellcome Trust – Medical Research Council Cambridge Stem Cell Institute. The authors would like to thank WTSI Cytometry Core Facility for their help with index cell sorting and the Core Sanger Web Team for hosting the cloud web application. The authors would also like to thank the CRUK Cambridge Institute Genomics Core Facility for their contribution in sequencing the data.

AUTHOR CONTRIBUTIONS

E.I.A. carried out the analysis; J.G.B. and L.F. performed experiments; H.A. generated a Hidden Markov model and DNN; P.L. oversaw implementation of the Hidden Markov model and DNN; J.G.B., E.I.A., and A.C. contributed to the discussion of the results and designed figures; A.C. conceived the study and wrote the manuscript. All authors approved the final version of the manuscript.

COMPETING FINANCIAL INTERESTS

The authors declare no competing financial interests.

REFERENCES

1. Orkin, S.H. & Zon, L.I. Hematopoiesis: An evolving paradigm for stem cell biology. *Cell* **132**, 631-644 (2008).
2. Kondo, M., Weissman, I.L. & Akashi, K. Identification of clonogenic common lymphoid progenitors in mouse bone marrow. *Cell* **91**, 661-672 (1997).
3. Akashi, K., Traver, D., Miyamoto, T. & Weissman, I.L. A clonogenic common myeloid progenitor that gives rise to all myeloid lineages. *Nature* **404**, 193-197 (2000).
4. Adolfsson, J. *et al.* Identification of Flt3(+) lympho-myeloid stem cells lacking erythromegakaryocytic potential: A revised road map for adult blood lineage commitment. *Cell* **121**, 295-306 (2005).
5. Mansson, R. *et al.* Molecular evidence for hierarchical transcriptional lineage priming in fetal and adult stem cells and multipotent progenitors. *Immunity* **26**, 407-419 (2007).
6. Notta, F. *et al.* Distinct routes of lineage development reshape the human blood hierarchy across ontogeny. *Science* **351**, 139-+ (2016).
7. Spangrude, G.J., Heimfeld, S. & Weissman, I.L. Purification and Characterization of Mouse Hematopoietic Stem-Cells. *Science* **241**, 58-62 (1988).
8. Guo, G.J. *et al.* Mapping Cellular Hierarchy by Single-Cell Analysis of the Cell Surface Repertoire. *Cell Stem Cell* **13**, 492-505 (2013).
9. Wilson, N.K. *et al.* Combined Single-Cell Functional and Gene Expression Analysis Resolves Heterogeneity within Stem Cell Populations. *Cell Stem Cell* **16**, 712-724 (2015).
10. Jaitin, D.A. *et al.* Massively Parallel Single-Cell RNA-Seq for Marker-Free Decomposition of Tissues into Cell Types. *Science* **343**, 776-779 (2014).
11. Paul, F. *et al.* Transcriptional Heterogeneity and Lineage Commitment in Myeloid Progenitors (vol 163, pg 1663, 2015). *Cell* **164**, 325-325 (2016).

12. Psaila, B. *et al.* Single-cell profiling of human megakaryocyte-erythroid progenitors identifies distinct megakaryocyte and erythroid differentiation pathways. *Genome Biol* **17** (2016).
13. Yamamoto, R. *et al.* Clonal Analysis Unveils Self-Renewing Lineage-Restricted Progenitors Generated Directly from Hematopoietic Stem Cells. *Cell* **154**, 1112-1126 (2013).
14. Treutlein, B. *et al.* Reconstructing lineage hierarchies of the distal lung epithelium using single-cell RNA-seq. *Nature* **509**, 371-375 (2014).
15. Trapnell, C. *et al.* The dynamics and regulators of cell fate decisions are revealed by pseudotemporal ordering of single cells. *Nat Biotechnol* **32**, 381-U251 (2014).
16. Macaulay, I.C. *et al.* Single-Cell RNA-Sequencing Reveals a Continuous Spectrum of Differentiation in Hematopoietic Cells. *Cell Rep* **14**, 966-977 (2016).
17. Leng, L. *et al.* MIF signal transduction initiated by binding to CD74. *J Exp Med* **197**, 1467-1476 (2003).
18. Shi, G.P. *et al.* Human Cathepsin-S - Chromosomal Localization, Gene Structure, and Tissue Distribution. *J Biol Chem* **269**, 11530-11536 (1994).
19. Wittamer, V., Bertrand, J.Y., Gutschow, P.W. & Traver, D. Characterization of the mononuclear phagocyte system in zebrafish. *Blood* **117**, 7126-7135 (2011).
20. Furze, R.C. & Rankin, S.M. Neutrophil mobilization and clearance in the bone marrow. *Immunology* **125**, 281-288 (2008).
21. Rosowski, E.E., Deng, Q., Keller, N.P. & Huttenlocher, A. Rac2 Functions in Both Neutrophils and Macrophages To Mediate Motility and Host Defense in Larval Zebrafish. *J Immunol* **197**, 4780-4790 (2016).
22. Kumar, S. *et al.* Cdc42 regulates neutrophil migration via crosstalk between WASp, CD11b, and microtubules. *Blood* **120**, 3563-3574 (2012).
23. Jones, R.A. *et al.* Modelling of human Wiskott-Aldrich syndrome protein mutants in zebrafish larvae using in vivo live imaging. *J Cell Sci* **126**, 4077-4084 (2013).
24. Grimm, T. *et al.* Dominant-negative Pes1 mutants inhibit ribosomal RNA processing and cell proliferation via incorporation into the PeBoW-complex. *Nucleic Acids Res* **34**, 3030-3043 (2006).
25. Brombin, A., Joly, J.S. & Jamen, F. New tricks for an old dog: ribosome biogenesis contributes to stem cell homeostasis. *Curr Opin Genet Dev* **34**, 61-70 (2015).
26. Ganis, J.J. *et al.* Zebrafish globin switching occurs in two developmental stages and is controlled by the LCR. *Dev Biol* **366**, 185-194 (2012).
27. Denker, B.M., Smith, B.L., Kuhajda, F.P. & Agre, P. Identification, Purification, and Partial Characterization of a Novel Mr 28,000 Integral Membrane-Protein from Erythrocytes and Renal Tubules. *J Biol Chem* **263**, 15634-15642 (1988).
28. Huang, H. & Cantor, A.B. Common Features of Megakaryocytes and Hematopoietic Stem Cells: What's the Connection? *J Cell Biochem* **107**, 857-864 (2009).
29. Renshaw, S.A. *et al.* A transgenic zebrafish model of neutrophilic inflammation. *Blood* **108**, 3976-3978 (2006).
30. Long, Q.M. *et al.* GATA-1 expression pattern can be recapitulated in living transgenic zebrafish using GFP reporter gene. *Development* **124**, 4105-4111 (1997).
31. Tamplin, O.J. *et al.* Hematopoietic Stem Cell Arrival Triggers Dynamic Remodeling of the Perivascular Niche. *Cell* **160**, 241-252 (2015).
32. Tang, Q. *et al.* Optimized cell transplantation using adult rag2 mutant zebrafish. *Nat Methods* **11**, 821-824 (2014).
33. Signer, R.A., Magee, J.A., Salic, A. & Morrison, S.J. Haematopoietic stem cells require a highly regulated protein synthesis rate. *Nature* **509**, 49-54 (2014).
34. Carroll, K.J. & North, T.E. Oceans of opportunity: Exploring vertebrate hematopoiesis in zebrafish. *Exp Hematol* **42**, 684-696 (2014).
35. Carmona, S.J. *et al.* Single-cell transcriptome analysis of fish immune cells provides insight into the evolution of vertebrate immune cell types. *Genome Res* (2017).

36. Hakes, L., Pinney, J.W., Lovell, S.C., Oliver, S.G. & Robertson, D.L. All duplicates are not equal: the difference between small-scale and genome duplication. *Genome Biol* **8** (2007).
37. Betancur, R.R. *et al.* The tree of life and a new classification of bony fishes. *PLoS Curr* **5** (2013).
38. Watkins, N.A. *et al.* A HaemAtlas: characterizing gene expression in differentiated human blood cells. *Blood* **113**, E1-E9 (2009).
39. Novershtern, N. *et al.* Densely Interconnected Transcriptional Circuits Control Cell States in Human Hematopoiesis. *Cell* **144**, 296-309 (2011).
40. Chen, L. *et al.* Transcriptional diversity during lineage commitment of human blood progenitors. *Science* **345**, 1580-+ (2014).
41. Nestorowa, S. *et al.* A single-cell resolution map of mouse hematopoietic stem and progenitor cell differentiation. *Blood* **128**, E20-E31 (2016).
42. Lin, H.F. *et al.* Analysis of thrombocyte development in CD41-GFP transgenic zebrafish. *Blood* **106**, 3803-3810 (2005).
43. Zhang, X.Y. & Rodaway, A.R. SCL-GFP transgenic zebrafish: in vivo imaging of blood and endothelial development and identification of the initial site of definitive hematopoiesis. *Dev Biol* **307**, 179-194 (2007).
44. Hall, C., Flores, M.V., Storm, T., Crosier, K. & Crosier, P. The zebrafish lysozyme C promoter drives myeloid-specific expression in transgenic fish. *Bmc Dev Biol* **7** (2007).
45. Walton, E.M., Cronan, M.R., Beerman, R.W. & Tobin, D.M. The Macrophage-Specific Promoter *mfap4* Allows Live, Long-Term Analysis of Macrophage Behavior during Mycobacterial Infection in Zebrafish. *Plos One* **10** (2015).
46. Dee, C.T. *et al.* CD4-Transgenic Zebrafish Reveal Tissue -Resident Th2-and Regulatory T Cell-like Populations and Diverse Mononuclear Phagocytes. *J Immunol* **197**, 3520-3530 (2016).
47. Bielczyk-Maczynska, E. *et al.* A loss of function screen of identified genome-wide association study Loci reveals new genes controlling hematopoiesis. *PLoS Genet* **10**, e1004450 (2014).
48. Picelli, S. *et al.* Smart-seq2 for sensitive full-length transcriptome profiling in single cells. *Nat Methods* **10**, 1096-1098 (2013).
49. Picelli, S. *et al.* Full-length RNA-seq from single cells using Smart-seq2. *Nat Protoc* **9**, 171-181 (2014).
50. Stachura, D.L. *et al.* Zebrafish kidney stromal cell lines support multilineage hematopoiesis. *Blood* **114**, 279-289 (2009).
51. Patro, R., Mount, S.M. & Kingsford, C. Sailfish enables alignment-free isoform quantification from RNA-seq reads using lightweight algorithms. *Nat Biotechnol* **32**, 462-U174 (2014).
52. Lun, A.T.L., Bach, K. & Marioni, J.C. Pooling across cells to normalize single-cell RNA sequencing data with many zero counts. *Genome Biol* **17** (2016).
53. Brennecke, P. *et al.* Accounting for technical noise in single-cell RNA-seq experiments. *Nat Methods* **10**, 1093-1095 (2013).
54. Buettner, F. *et al.* Computational analysis of cell-to-cell heterogeneity in single-cell RNA-sequencing data reveals hidden subpopulations of cells. *Nat Biotechnol* **33**, 155-160 (2015).
55. Haghverdi, L., Buettner, F. & Theis, F.J. Diffusion maps for high-dimensional single-cell analysis of differentiation data. *Bioinformatics* **31**, 2989-2998 (2014).
56. Reimand, J. *et al.* g:Profiler-a web server for functional interpretation of gene lists (2016 update). *Nucleic Acids Res* **44**, W83-89 (2016).

FIGURE LEGENDS

Figure 1. Pseudotime ordering reveals a gradual transition of cells from immature to more differentiated within the myeloid branch

a) Experimental strategy for sorting single cells from transgenic zebrafish lines. Cells were harvested from a single kidney of each line and sorted for expression of the fluorescent transgene. Index sorting was used to dispense single cells into a 96 well plate and these were subsequently processed for RNA-seq analyses. b) Five cell states were predicted using the Monocle2 algorithm for temporal analyses of single cell transcriptomes. c) Analysis of genes that are differentially expressed across the five states (given the same colour code used in b) reveals GO terms (inner circle) that are highly pertinent to specific cell types. The outer circle shows examples of May-Grünwald Giemsa stained cells from kidneys of transgenic lines that largely label each particular cell type. d) Jitter plots showing the expression (y axis) of differentially expressed marker genes in each cell type (x axis). Each dot in the jitter plot shows the expression of the gene $\log_{10}(\text{counts} + 1)$ in each cell.

Figure 2. The distribution of cells from different transgenic lines modelled by Monocle

a) The trajectories of cell states predicted by Monocle are shown in grey for each transgenic line used, with the associated cell types labelled in blue. The percentage of cells from each transgenic line contributing to each state is given next to the relevant trajectory. b) Pie charts showing the contribution of transgenic lines to each cell type. The colour code relates to the colours given in the headers for each transgenic line used in (a).

Figure 3. Cells within distinct states have different repopulation potentials

a) Experimental strategy for the adult transplantation experiment. Kidneys were dissected from transgenic donor fish and sorted for cells expressing the fluorescent transgene. Positive cells were collected and injected into sub-lethally irradiated *rag2^{E450fs/-}* fish. b) Assessment for engraftment was made one day, four- and 14 weeks post transplantation using flow cytometry. Successfully engrafted fluorescent donor cells were isolated at four weeks PI by index sorting single cells into a microtitre plate for subsequent RNA-seq analyses. c) Distribution of *runx1+* cells, from non-transplanted (right) and transplanted (left) fish, modelled by Monocle.

Figure 4. Transcriptionally similar cells display different probabilities of transitioning to another state. The approximate positions of the cell states identified by Monocle 2 are shown in the insert. The graph shows distribution of HSPCs in pseudotime. Cells are coloured based on their probability to transition to a specific cell type. Cells that

have less than 75% probability to transition to any given cell type are defined as “unspecified”. The proportion of cells that belongs to each of the predicted states is shown in the Pie chart.

Figure 5. Lineage differentiation is defined by two main trends in gene expression

a) Heatmap of genes whose expression changed dynamically during pseudotime in each of the four branches. b) Graph showing the average expression pattern of the dynamically expressed genes that follow the same trend across pseudotime. For each of the cell states, one gene is presented that follows one of the two main trends. Standard error is shown as a gray area around the trend lines. c) Heatmap of expression of 168 genes annotated as “ribosomal proteins” genes in pseudotime in each of the four branches.

Figure 6. Conservation analysis of zebrafish genes differentially expressed in the main blood cell types.

a) Percentage of zebrafish protein-coding genes (specific for distinct blood cell types, as well as non-differentially expressed) with orthologs in other vertebrate species. b) The total number of paralogs duplicated exclusively pre- (green) and post ray-finned speciation (red). The numbers 1-7 mark the number of cell types (erythrocytes, monocytes, neutrophils, thrombocytes, HSPCs, T-cells and NK cells) in which the duplicated genes are expressed. c) The percentage of conserved vs diverged genes duplicated exclusively post speciation (fish specific genes).

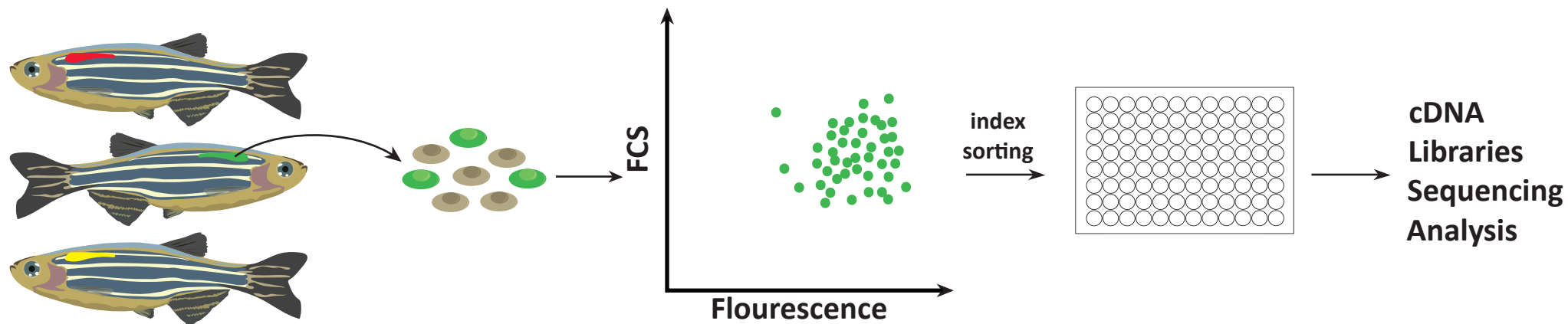
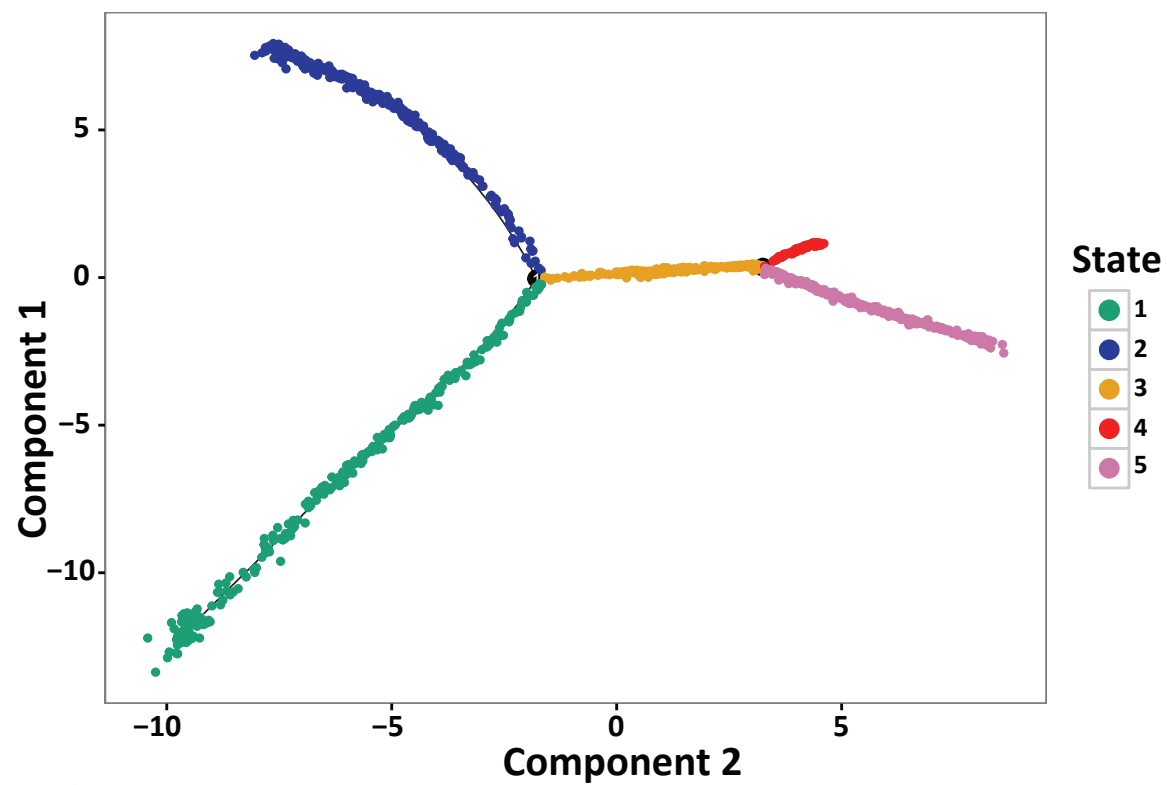
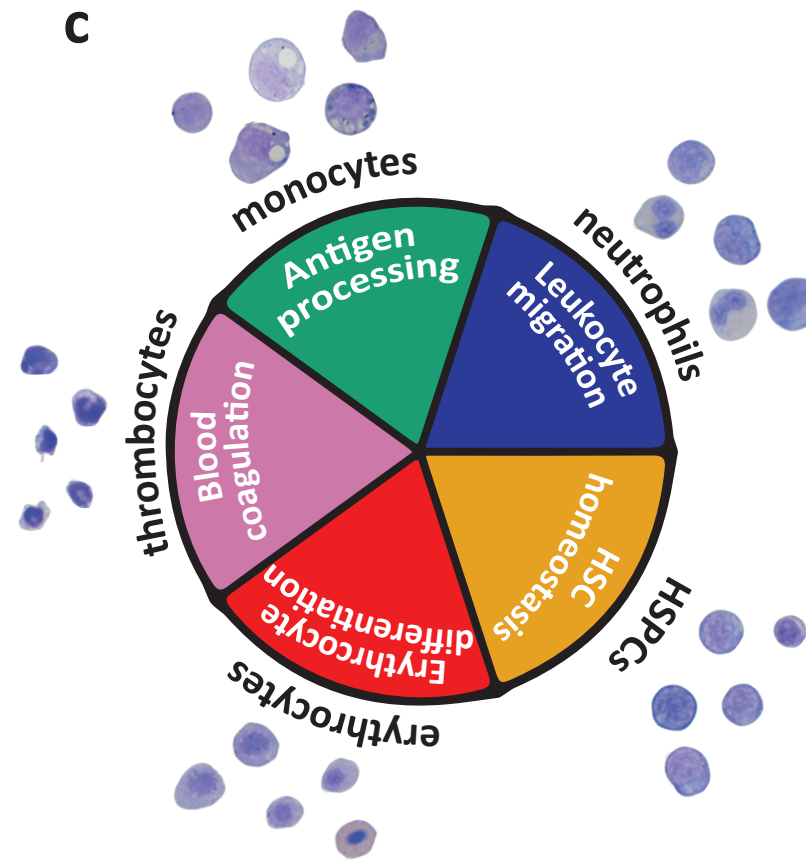
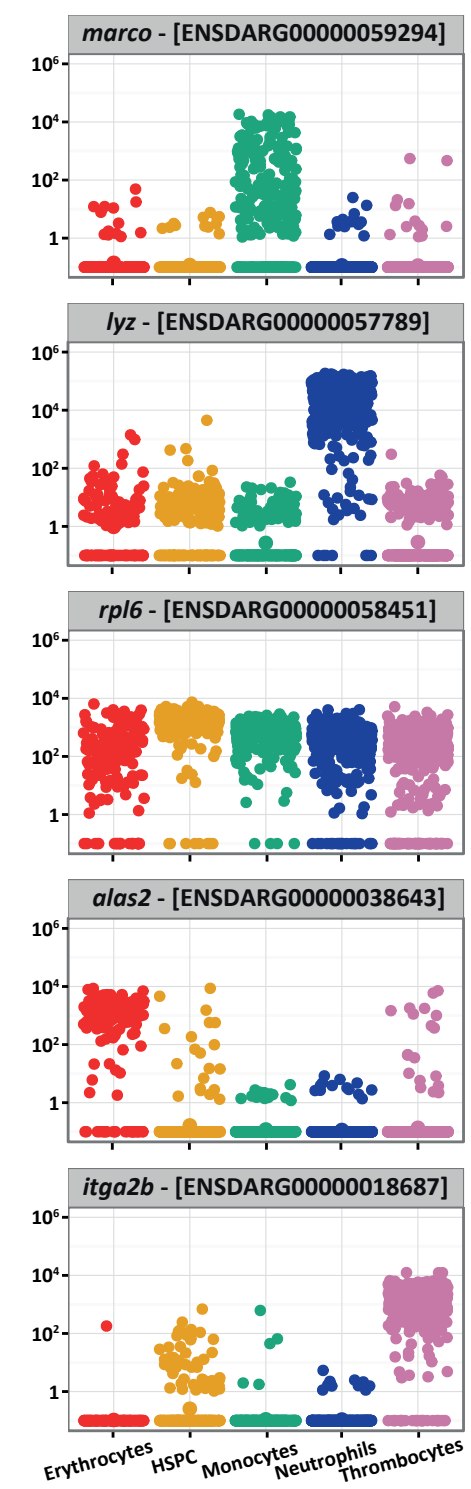
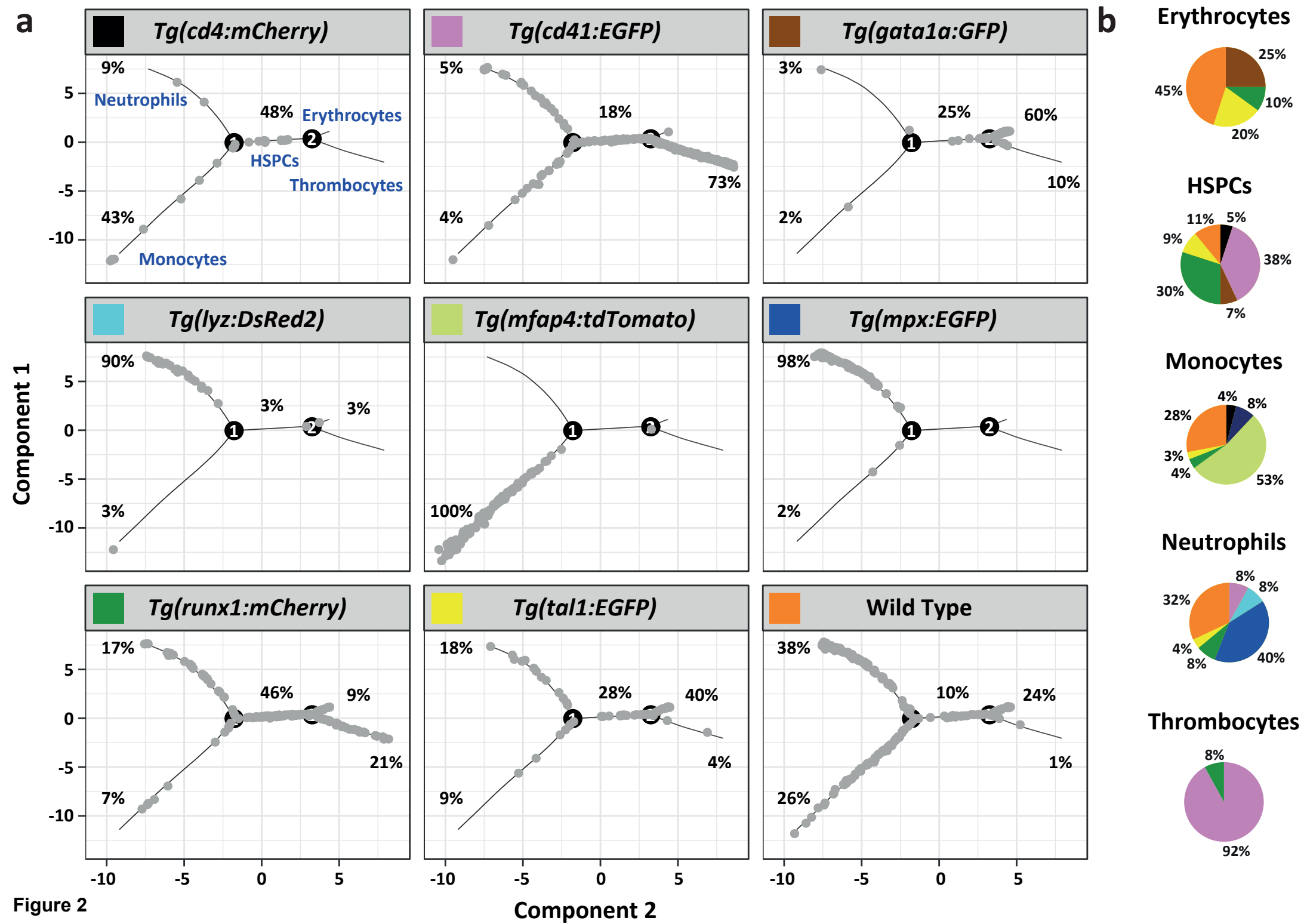
a**b**

Figure 1

c**d**



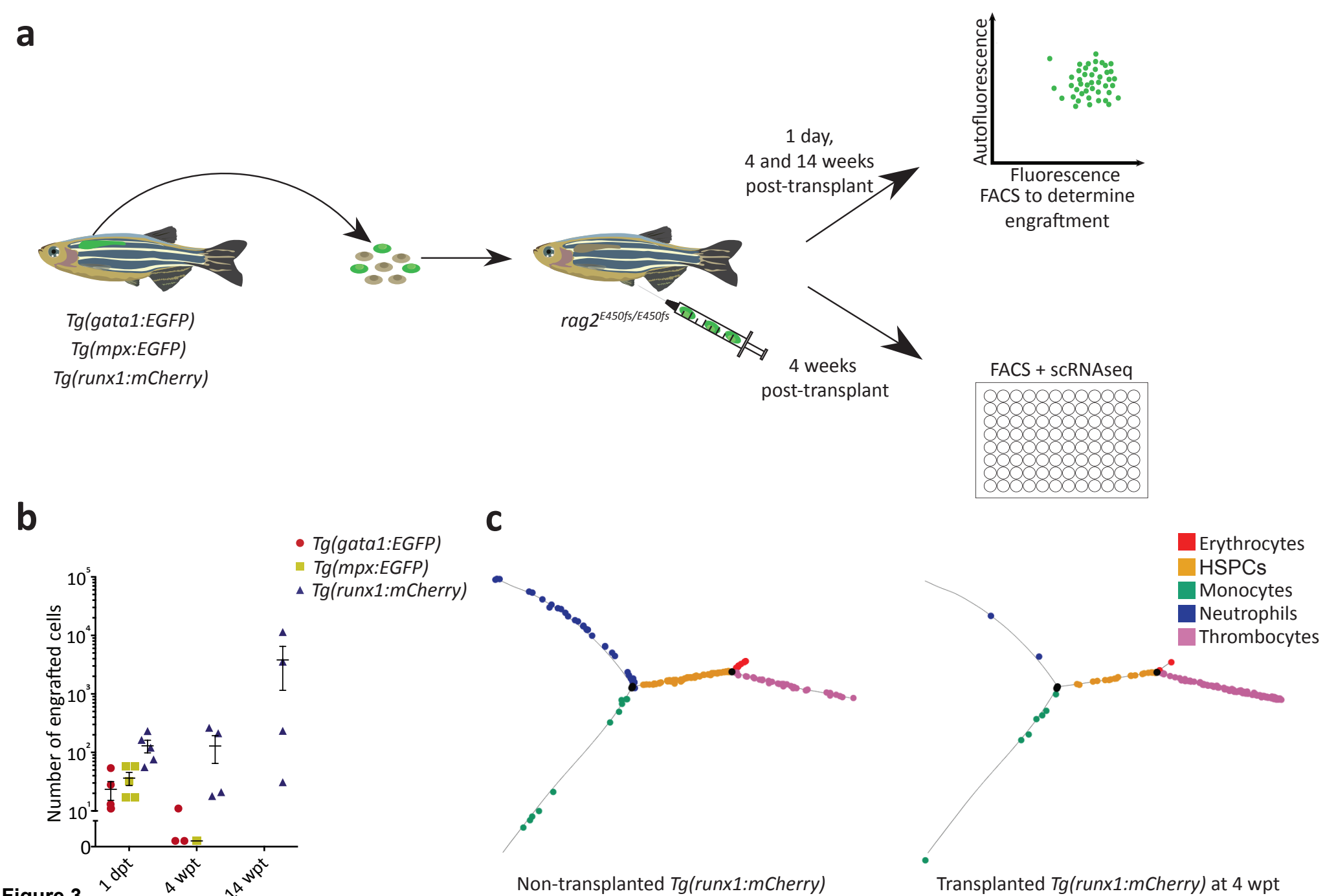
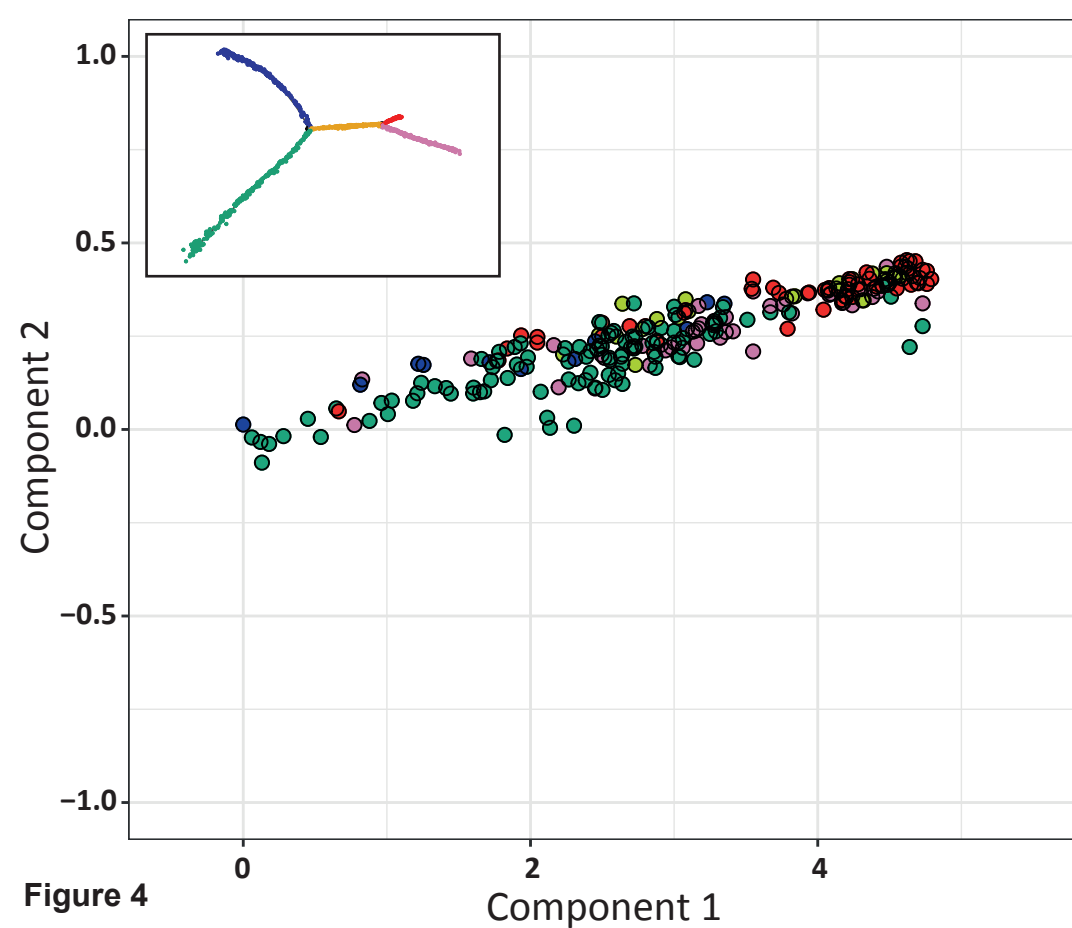
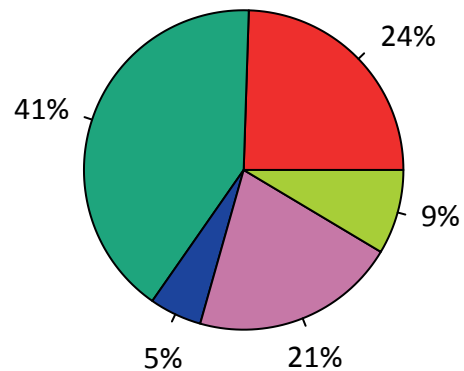


Figure 3



Predicted State

- Erythrocytes
- Monocytes
- Neutrophils
- Thrombocytes
- Unspecified



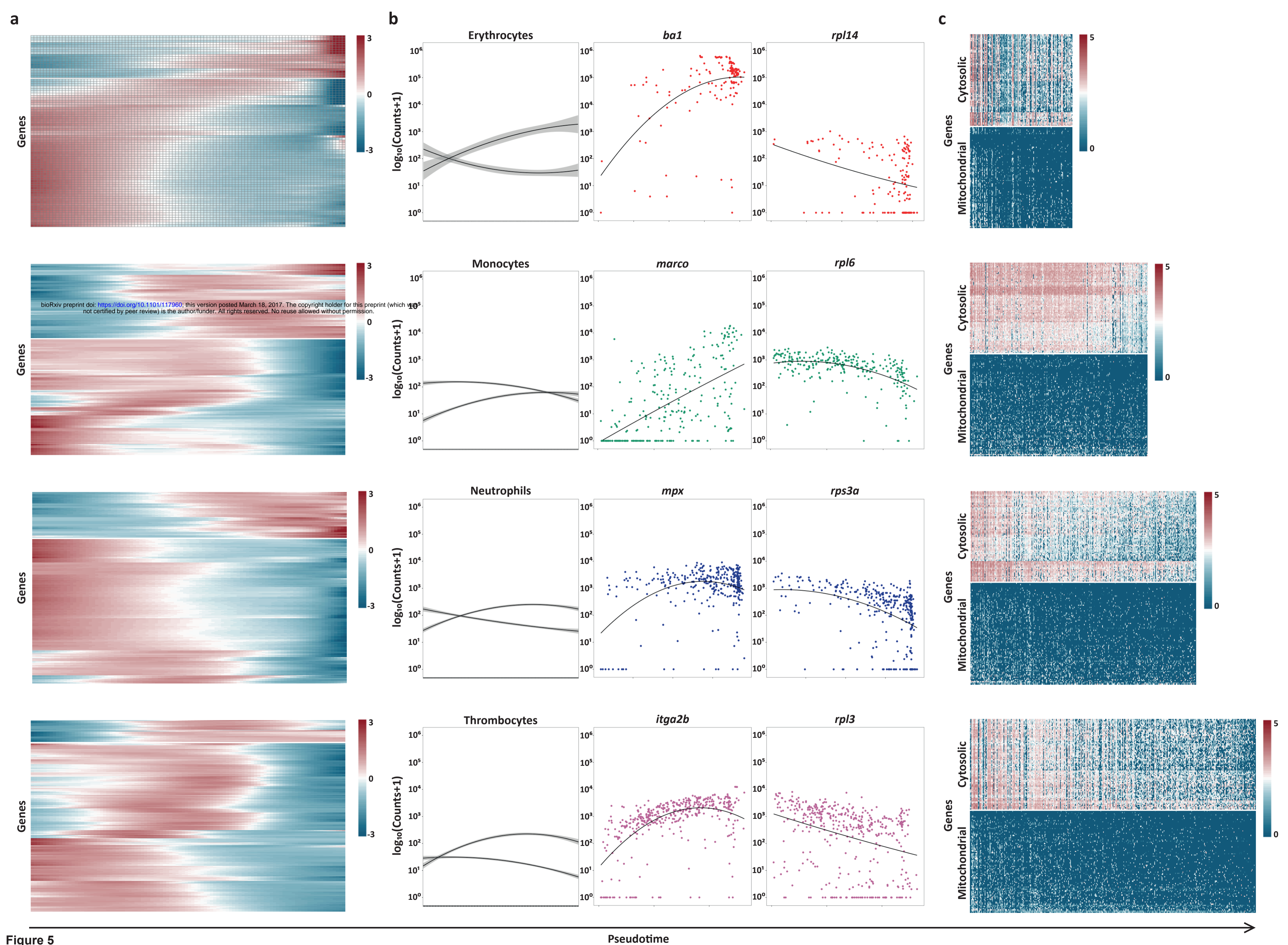


Figure 5

a

Proportion of Genes with Orthologues (%)

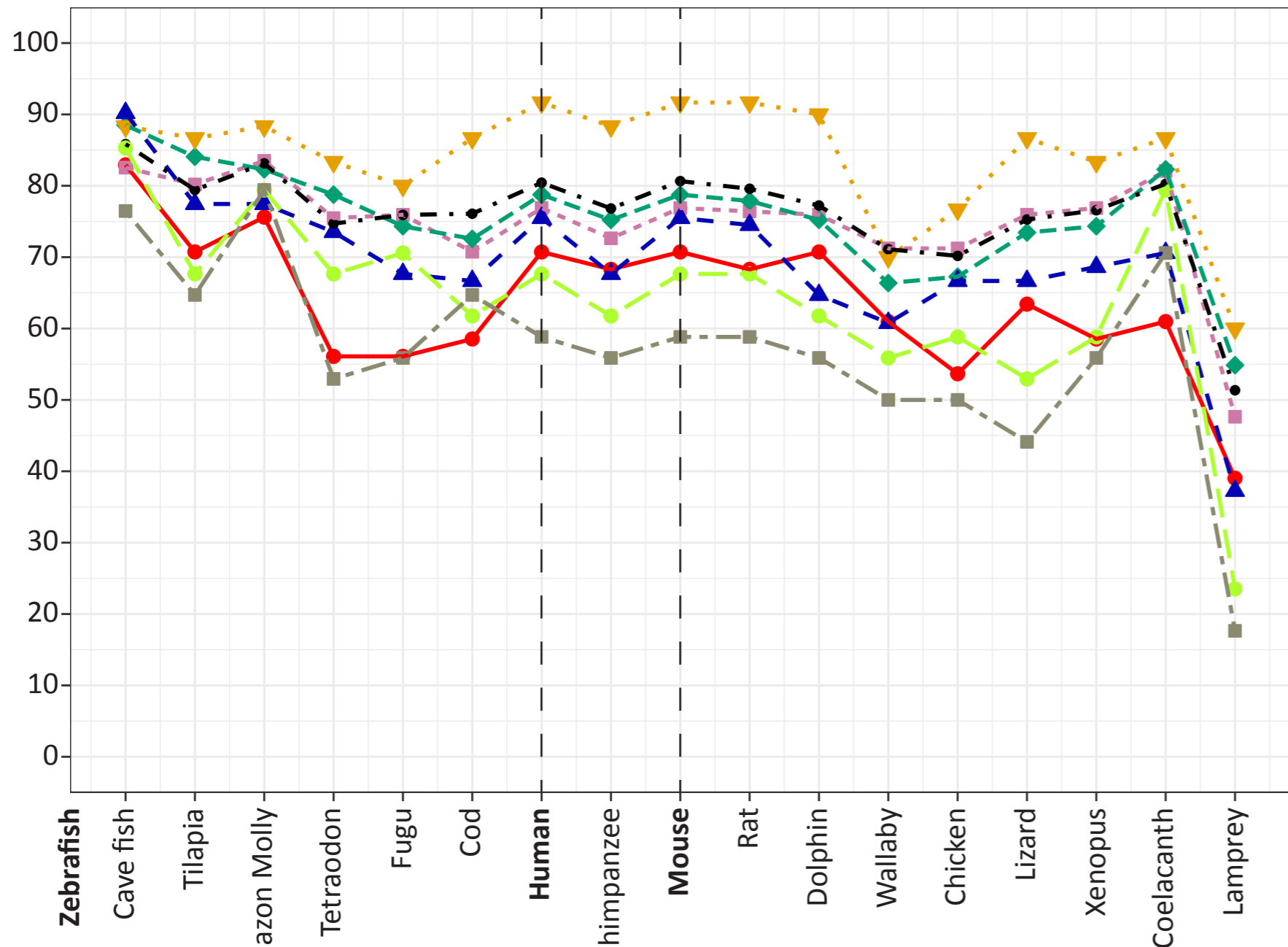
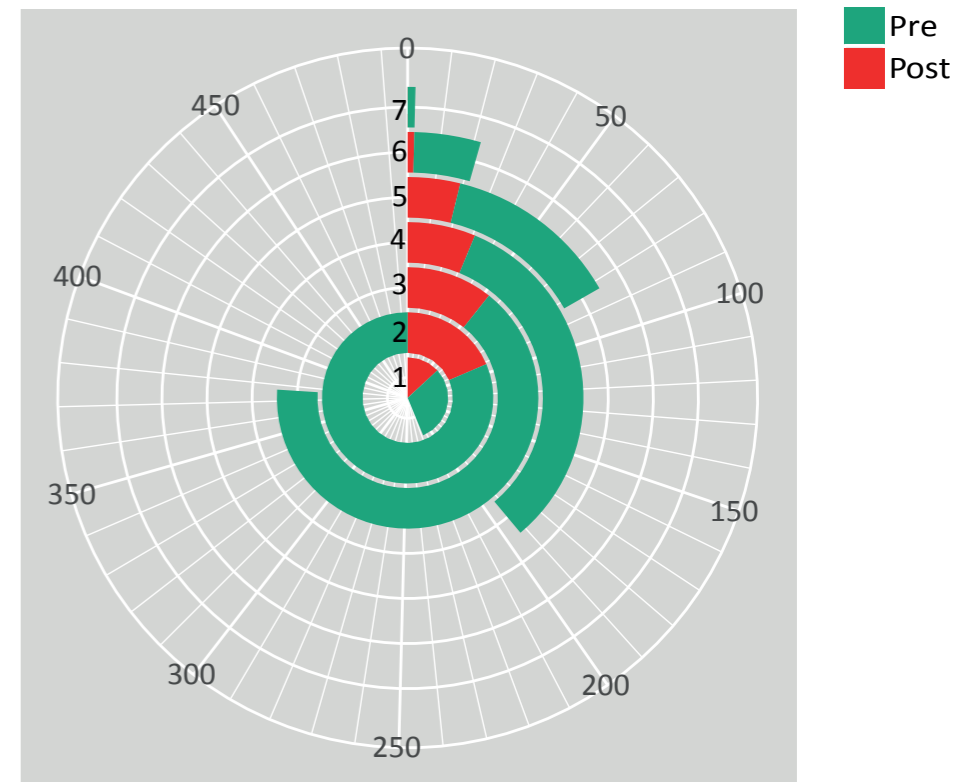


Figure 6

b



c

