



# MIT Open Access Articles

## *A systems biology pipeline identifies regulatory networks for stem cell engineering*

The MIT Faculty has made this article openly available. **Please share** how this access benefits you. Your story matters.

<b>Citation</b>	Kinney, Melissa A. et al. "A systems biology pipeline identifies regulatory networks for stem cell engineering." Nature Biotechnology 37, 7 (July 2019): 810–818 © 2019 Springer Nature
<b>As Published</b>	<a href="http://dx.doi.org/10.1038/s41587-019-0159-2">http://dx.doi.org/10.1038/s41587-019-0159-2</a>
<b>Publisher</b>	Springer Science and Business Media LLC
<b>Version</b>	Author's final manuscript
<b>Citable link</b>	<a href="https://hdl.handle.net/1721.1/125922">https://hdl.handle.net/1721.1/125922</a>
<b>Terms of Use</b>	Article is made available in accordance with the publisher's policy and may be subject to US copyright law. Please refer to the publisher's site for terms of use.



Published in final edited form as:

*Nat Biotechnol.* 2019 July ; 37(7): 810–818. doi:10.1038/s41587-019-0159-2.

## A systems biology pipeline identifies regulatory networks for stem cell engineering

Melissa A. Kinney<sup>1,2,3</sup>, Linda T. Vo<sup>1,2,4</sup>, Jenna M. Frame<sup>1,2,5</sup>, Jessica Barragan<sup>1,2</sup>, Ashlee J. Conway<sup>1,2</sup>, Shuai Li<sup>6,7</sup>, Kwok-Kin Wong<sup>6,7</sup>, James J. Collins<sup>8,9,10,11</sup>, Patrick Cahan<sup>12</sup>, Trista E. North<sup>1,2,5</sup>, Douglas A. Lauffenburger<sup>#,3</sup>, George Q. Daley<sup>#,1,2,4,13</sup>

<sup>1</sup>Stem Cell Program, Boston Children's Hospital, Boston, MA

<sup>2</sup>Division of Hematology/Oncology, Boston Children's Hospital and Dana Farber Cancer Institute, Boston, MA

<sup>3</sup>Department of Biological Engineering, Massachusetts Institute of Technology, Cambridge, MA

<sup>4</sup>Harvard Medical School, Boston, MA

<sup>5</sup>Department of Pathology, Beth Israel-Deaconess Medical Center, Boston MA

<sup>6</sup>Department of Medical Oncology, Dana-Farber Cancer Institute, Boston, MA

<sup>7</sup>Laura & Isaac Perlmutter Cancer Center, NYU Langone Medical Center, New York, NY.

<sup>8</sup>Institute for Medical Engineering & Science, Department of Biological Engineering, Synthetic Biology Center, Massachusetts Institute of Technology

<sup>9</sup>Harvard-MIT Program in Health Sciences and Technology

<sup>10</sup>Broad Institute of MIT and Harvard

<sup>11</sup>Wyss Institute for Biologically Inspired Engineering, Harvard University

<sup>12</sup>Department of Biomedical Engineering, Johns Hopkins University, Baltimore, MD

<sup>13</sup>Harvard Stem Cell Institute, Boston, MA

### Abstract

A major challenge for stem cell engineering is achieving a holistic understanding of the molecular networks and biological processes governing cell differentiation. To address this challenge, we describe a computational approach that combines gene expression analysis, prior knowledge from proteomic pathway informatics, and cell signaling models to delineate key transitional states of differentiating cells at high resolution. Our network models connect sparse gene signatures with corresponding, yet disparate, biological processes to uncover molecular mechanisms governing cell fate transitions. This approach builds upon our earlier CellNet and recent trajectory-defining

<sup>#</sup>Correspondence should be addressed to: lauffen@mit.edu or George.Daley@childrens.harvard.edu.

**Contributions.** M.A.K., D.A.L., and G.Q.D. conceived the project. M.A.K., L.T.V., J.M.F., J.B., A.J.C., and S.L., performed experimental work and data interpretation. K-K.W., J.J.C., P.C., T.E.N., D.A.L., and G.Q.D. supervised research and participated in project planning. M.A.K., T.E.N., D.A.L., and G.Q.D. prepared the manuscript.

**Data Availability.** All RNA-seq data has been deposited to the Gene Expression Omnibus (GEO) database under GSE108128.

**Competing financial interests.** The authors declare no competing financial interests.

algorithms, as illustrated by our analysis of hematopoietic specification along the erythroid lineage, which reveals a role for the EGF receptor family member, ErbB4, as an important mediator of blood development. We experimentally validate this prediction and perturb the pathway to improve erythroid maturation from human pluripotent stem cells. These results exploit an integrative systems perspective to identify new regulatory processes and nodes useful in cell engineering.

---

Stem cell biology, cell engineering, and regenerative medicine often invoke developmental principles to differentiate cells toward target identities. However, much remains to be learned about how signaling pathways integrate to determine cell fate<sup>1</sup>. The past decade of cell engineering has shown that expression of individual genes, or sets of genes, is often insufficient to functionally reprogram cell identity<sup>2,3</sup>, underscoring the need for new approaches to quantitatively describe and manipulate cell state. We previously established CellNet<sup>4-6</sup> to assess the fidelity of engineered cells by interrogating key gene regulatory networks (GRNs) that define native populations. CellNet extracts cell type-specific GRNs from transcriptional profiling data, compares the GRNs to those of bona fide primary cells and tissues to assign a similarity metric, and identifies dysregulated transcriptional regulators that account for the differences between engineered cells and their native counterparts. The network-level CellNet algorithm confers robustness to biological and technical variability and encodes topological information about regulator-target relationships. A limitation of CellNet is that training data consisting of a small number of terminal cell and tissue types obscures the phenotypic heterogeneity that arises during dynamic biological processes like cell differentiation. More recent efforts have aimed to describe intermediate developmental states using trajectory-based methods, which employ cell-cell similarity metrics to infer dynamics<sup>7-10</sup>. However, these algorithms rely on single-cell transcriptomics to provide sufficiently powered datasets and largely forgo network analytics.

Here we extend CellNet to quantitatively define network dynamics along a differentiation pathway. We show that publicly accessible gene expression datasets capture population-level differentiation states with high dynamic resolution and broad biological scope, including responses across a spectrum of experimental variables like chemical and genetic perturbations. Our pipeline goes beyond the establishment of GRNs to enable quantification of differentiation dynamics and identification of key signaling pathways governing cell fate changes. We apply this otherwise general approach to characterize erythropoiesis, a dynamic process that generates red blood cells (RBCs) throughout the lifetime of the organism. We focused on this system because its temporal stages of differentiation, defined by distinct immunophenotypes, have been comprehensively characterized<sup>11</sup>. Our analyses confirm key processes involved in distinct stages of erythropoiesis and elucidate novel dynamic patterns of gene expression. To improve erythroid maturation *in vitro*, we constructed an interaction network connecting the dynamic molecular signatures that stratify late erythroblasts from reticulocytes. Our network analytics identifies a role for ErbB signaling during erythropoiesis, which we validate in human, murine and zebrafish models and apply to iPS-derived RBC maturation.

## RESULTS

### CellNet delineates stem cells and progeny

To analyze the dynamics of stem cell differentiation, we began by establishing GRNs for hematopoietic stem cells and differentiated progeny using CellNet as previously described<sup>4–6</sup>. We augmented the original CellNet compendium<sup>4</sup> of microarray datasets from 16 human cell and tissue types to include 164 publicly available erythroid microarray datasets (Supplementary Tables 1–2). The corresponding data represented a broad phenotypic range, including those manually classified by the discrete designations of early (CFU-E), intermediate (IntE) and late (LateE) erythroid progenitors, as well as reticulocytes and the K526 erythroleukemia cell line (Supplementary Table 3). Our rationale for augmenting the original compendium arises from the paucity and lack of biological variability in publicly available erythroid-specific sequencing datasets at the time of data compilation and manuscript submission.

Application of the original CellNet classifier identified erythroid cells as hematopoietic stem and progenitor cells (HSPCs) with high probability (Supplementary Fig. 1a). However, after re-training the classifier with the augmented compendium and establishing an erythroid-specific GRN, CellNet robustly distinguished HSPCs and erythroid cell types (Supplementary Fig. 1b,c), with little overlap between the two GRNs other than cofactors mediating the canonical “GATA switch”<sup>12</sup> that governs erythroid specification from HSPCs (Supplementary Fig. 1d–e,g). The erythroid GRN comprised 235 genes that were highly enriched for biological processes such as hemoglobin synthesis, oxygen transport, cell cycle, and hematopoietic development (Supplementary Fig. 1f), with subnetworks governed by 17 transcription factors, including canonical erythroid factors<sup>13</sup> such as GATA1, LMO2, TAL1, and KLF1. We also identified several factors not fully characterized in erythropoiesis, including HES6<sup>14</sup>, CDT1<sup>15</sup>, and SREBF1<sup>16</sup>. These data demonstrate that the CellNet algorithm can be readily augmented by the addition of new classifiers and accurately distinguishes biologically relevant GRNs, even among closely related cell types.

### GRNs capture cell state dynamics

The augmented CellNet algorithm classified all stages of erythropoiesis with high probability<sup>4</sup>, (Supplementary Fig. 1h–j). We hypothesized that, analogous to the extraction of cell type-specific GRNs, we could identify subnetworks, or smaller gene modules, within the erythroid GRN that would enable us to classify distinct stages of differentiation. To dissect the erythroid GRN, we projected the data into a Principal Component space (Fig. 1b). PC1 (35.2% variance) was highly correlated with the GSEA hallmark pathway for hemoglobin metabolism (Pearson’s  $r = 0.91$ ; Fig. 1b, Supplementary Fig. 2a), and PC2 (19.5% variance) correlated with MYC targets (Pearson’s  $r = -0.8$ ; Fig. 1b, Supplementary Fig. 2b). Unsupervised Gaussian Mixture Model (GMM)–based clustering identified six discrete phenotypes, which were significantly enriched for manual, literature-based designations of erythroid stage—C1–C2: CFU-E/early proerythroblast; C3–C4: intermediate proerythroblast; C5: late erythroblast; C6: reticulocyte (Fig. 1c, Supplementary Table 3). K562, an erythroleukemia cell line, clustered in C1, and studies of hemoglobin-perturbed cells clustered in C3, a PC1-shifted intermediate cluster. We therefore focused on clusters

C2, C4, C5 and C6 for the purpose of studying physiological and developmental erythropoiesis.

In contrast to whole genome–based dimensionality reduction techniques commonly used in trajectory algorithms, the erythroid GRN served as a feature selection upstream of PCA, which identifies genes correlated with developmental stages. By GMM clustering, erythroid network genes clustered into 3 distinct groups (Supplementary Table 4), with early (G1) and intermediate (G2) differentiation clusters associated with cell cycle and hemoglobin synthesis, respectively (Fig. 1d, Supplementary Table 5). The reticulocyte cluster (G3) comprised genes that were not significantly enriched for any biological processes. Likewise, ranking gene importance to each phenotypic cluster (Supplementary Fig. 2d–g; Supplementary Table 6) failed to yield annotations for the reticulocyte cluster (C6). We also implemented K-means clustering to identify sets of genes with similar dynamic expression across biological clusters (C2, C4, C5 and C6) and identified coordinated regulation of genes related to processes such as stress responses, autophagy, and apoptosis during differentiation (Supplementary Fig. 2h).

We confirmed that dimensionality reduction similarly captured biologically meaningful clusters in developmentally staged, purified populations analyzed by bulk RNA-seq (GSE53983). Sample localization in the PCA space was driven by similar genes, as shown by significant correlations with biological analogues from microarray data—proerythroblast: C2/S1/S2; intermediate erythroblast: C4/S3; late erythroblast: C5/S4; positive Pearson's  $r$ , with  $p < 0.05$  (Supplementary Fig. 3C). Of note, expert knowledge was required for the interpretation and comparison of biologically analogous samples across microarray and RNA-seq datasets (e.g. C4 and C5 both correlate with S3; C2 correlates with both S1 and S2), likely because the granularity and variance are strongly tied to the data source and experimental design. This further highlights the need for large data compendia, including purified populations, primary cells, in vitro–differentiated cells, genetically perturbed cells, and rare/unique populations (i.e. reticulocytes), to fully sample the biological space within a given cell type. We also demonstrated that, in addition to compatibility with different data types, our pipeline is generally applicable to other biological systems (Supplementary Fig. 4, Supplementary Tables 7–9).

To derive dynamic network models, we exploited topological regulatory information encoded in the erythroid GRN as a secondary layer atop the lineage-correlated loadings (Fig. 1e). Early erythropoiesis is dominated by a few, highly connected regulators, with more distributed regulation during the proerythroblast (C4) and late erythroblast (C5) stages (Fig. 1e). Canonical regulators, such as GATA1 and E2F2, are also highly connected in the mature (C6) network, suggesting that early regulators impart persistent influence (Fig. 1f). Indeed, GATA1 has been implicated in erythroid maturation, with a distinct network from that of early erythroid specification<sup>17</sup>. E2F2 also recurs as a central regulator in both the intermediate erythroblast (C4) and reticulocyte (C6) stages; however, there is a clear re-wiring of its targets from a diffuse cluster of co-regulated genes to a more compact network during maturation (Fig. 1f). Thus, the integration of network biology and GRN-based feature selection with dimensionality reduction uncovers dynamic changes in network activity and architecture accompanying cell fate changes.

## Identification of pathways mediating cell fate transitions

We further explored the capacity for network analytics to identify biological processes that mediate stem cell differentiation. We focused on the late erythroblast (C5) to reticulocyte (C6) transition, as relatively little is known about the integrated mechanisms controlling terminal erythroid maturation. Moreover, microarray datasets derived from the *in vivo* reticulocyte transcriptome<sup>18</sup> (Supplementary Fig. 5) provided comparisons that are not readily accessible, owing to such transient and mobile populations.

To construct signatures of this transition, we employed the Least Absolute Shrinkage and Selection Operator (LASSO), as a feature selection method that minimizes covariate correlation. The resulting 27 gene signature (Fig. 2a; Supplementary Table 10) accurately predicted the late erythroblast and reticulocyte cell states without overfitting, based on a Partial Least Squares Discriminant (PLSDA) model (Supplementary Fig. 6). This method produced a sparse gene set that lacked unifying annotations. We therefore adopted a 'bottom-up' approach using local network information to connect our signature genes (Fig. 2b). This propagation of LASSO targets is similar to network biology approaches to predict drug targets and disease-associated genes and is based on the hypothesis that genes in close proximity topologically are functionally related<sup>19–21</sup>.

To identify common regulators, we investigated the local topology of the first-order subnetwork in the global CellNet GRN, from which cell type specific GRNs were originally identified (Fig. 2c). Contrary to our hypothesis, there was largely a one-to-one connection between all connected regulators and LASSO targets (Fig. 2d). The few statistically enriched genes belonged to networks of co-regulated transcription factors, such as the pluripotency factors (i.e. NANOG, SOX2, LIN28)<sup>22,23</sup> which are associated with a single LASSO target, SALL2. This topology suggests that LASSO targets are associated with discrete biological processes, rather than being downstream of common regulators. Accordingly, this same analysis identified common regulators between ontologically related genes (Supplementary Fig. 7). However, further dissection of gene modules with modest co-regulation of LASSO targets revealed that late erythroblast targets (Fig 2e; Module 1) were associated with regulators of hematopoietic differentiation and P53-apoptotic pathways, whereas reticulocyte LASSO genes (Fig. 2e; Module 2) were downstream of metabolic and lipid pathways important for RBC maintenance.

Based on this largely one-to-one topology of the transcriptional regulator-target network, we hypothesized that common signaling networks may lie upstream of the transcriptional layer. Therefore, we generated an interaction network using the STRING database (Fig. 2f, Supplementary Fig. 8, Supplementary Table 11). We employed the Prize Collecting Steiner Forest (PCSF) algorithm<sup>24</sup>, which is particularly suited for modeling multiple, independent pathways acting in synergy toward a unified biological response. The resulting network was enriched for biological processes such as apoptotic signaling, stress responses, and cell cycle, consistent with prior analyses (Fig. 2g). Uniting these processes, P53 is a highly interconnected central node ( $p < 0.001$ ; Supplementary Fig. 8f). The network was also significantly enriched for Reactome signaling pathways relevant for erythropoiesis, such as Notch<sup>25</sup>, Rho<sup>26</sup>, TGF $\beta$ <sup>27</sup> and BCR<sup>28</sup>, as well as novel candidate pathways, including EGFR/ErbB4, TLR and RIG-I/MDA5 (Fig. 2h).

Finally, we used a ‘guilt by association’ approach to define networks that were highly correlated with the LASSO signature (Fig. 2i, Supplementary Fig. 9, Supplementary Table 12). Highly enriched transcription factor binding (ENCODE and ChEA) and kinase regulation (LINCS L1000) further implicated proliferative and apoptotic processes (i.e. E2F, P53 and FOXM1/WEE1). Moreover, several enriched kinases included members of the MAPK/ERK pathway (i.e. SRC, ErbB3/ErbB4), and the ligand activation signatures (EGF, TGFA, BTC) further supported a role for ErbB signaling in regulation of the coexpression network (Fig. 2j). This analysis demonstrates the utility of combining sparse gene signatures with network propagation approaches to identify novel biological processes that potentially mediate dynamic fate changes, hence establishing hypotheses to be experimentally confirmed.

### **ErbB4 is necessary for efficient erythropoiesis**

Although our network models identified several enriched signaling nodes and candidate pathways in erythroid maturation, the preponderance of evidence pointed to ErbB signaling. Although significantly enriched, EGFR/ErbB4 was not among the top candidate pathways (Fig. 2h; comprehensive list in Supplementary Table 8); however, when combined with expert knowledge that ErbB signaling is frequently associated with P53<sup>29</sup>, and the apoptotic<sup>30</sup> and proliferative<sup>31</sup> processes that were repeatedly identified in our network models, ErbB signaling emerged as a lead candidate. To determine whether ErbB signaling was necessary for erythroid maturation, we perturbed erythroblasts differentiated from bone marrow HSPCs (CD34<sup>+</sup>) with ErbB inhibitors (Fig. 3a). Maturation (GlyA+CD71<sup>-</sup>) was only affected by pan-ErbB inhibitors (Afatinib, Dacomitinib, Neratinib), implicating ErbB4 rather than EGFR/ErbB2.

ErbB4 signaling has not previously been implicated in blood development or homeostasis. We characterized ErbB4 in human, mouse, and zebrafish erythropoiesis. Using an erythroid in vitro differentiation protocol for human HSPCs<sup>32</sup>, we observed increasing ERBB4 mRNA expression as erythroid cells matured (Supplementary Fig. 10a). Native human bone marrow erythroid fractions also exhibited increased ERBB4 expression in the most mature population (GlyA<sup>+</sup>CD71<sup>-</sup>) (Supplementary Fig. 10b). Reciprocally, pharmacological inhibition of ErbB signaling with Neratinib for one week shifted the bone marrow differentiation profile in mice, with an increase in immature and a decrease in mature erythroid populations (Fig. 3b), as well as changes to the peripheral hematopoietic fractions (Supplementary Fig. 11).

We next determined whether ErbB4 signaling also functioned during erythroid ontogeny, a process which initiates in multiple waves from restricted progenitors during embryogenesis<sup>33</sup>. Morpholino inhibition of ErbB4 in zebrafish embryos significantly decreased the frequency of Gata1<sup>+</sup> erythroid cells (Supplementary Fig. 12a) and of more differentiated globin-expressing cells (Fig. 3c–d) at 48–56 hours post fertilization (hpf), without affecting neutrophils (Supplementary Fig. 12b). These data indicate that ErbB4 signaling is necessary for robust erythropoiesis during embryonic and adult hematopoiesis.

### ErbB4 deficiency induces stress erythropoiesis

To more stringently characterize ErbB4 in adult erythropoiesis, we employed a genetic mouse model derived via  $\alpha$ MHC-driven expression of human HER4 to circumvent embryonic lethality from heart defects in the whole body ErbB4 knockout (ErbB4<sup>-/-</sup>-HER4<sup>heart</sup>)<sup>34</sup>. Consistent with the effects of Neratinib treatment, we observed an increase in early proerythroblast populations, with fewer mature orthochromatic and normoblastic cells in the ErbB4<sup>-/-</sup> bone marrow (Supplementary Fig. 13a,b). Nucleated RBCs and a high percentage of reticulocytes were present in peripheral blood, indicating moderate stress erythropoiesis in homozygotes, and blood counts revealed significant changes in hemoglobin distribution (Fig. 4a, Supplementary Fig. 13c). ErbB4<sup>-/-</sup> mice had enlarged spleens (Supplementary Fig. 13d), a >2-fold expansion of early erythroblasts (GlyA+CD71+; gate II) (Fig. 4b) and overcrowded red pulp (Fig. 4c), suggesting extramedullary erythropoiesis. Morphological analysis demonstrated early developmental blocks across multiple lineages in ErbB4<sup>-/-</sup> bone marrow (Fig 4d). CD41<sup>+</sup> megakaryocytes in ErbB4-deficient spleen ( $p = 0.007$  compared to wild type; Fig. 4e) decreased significantly, accompanied by a myeloid-skewed leukocyte profile and increased platelets in the periphery (Fig. 4f). These results demonstrate dysregulated multi-lineage hematopoietic phenotypes in ErbB4<sup>-/-</sup>-HER4<sup>heart</sup> mice.

### Mitotic and proliferative processes downstream of ErbB matures iRBCs

To interrogate the molecular mechanisms downstream of ErbB, we performed global gene expression analysis of *in vitro* differentiated RBCs perturbed with pan versus selective inhibitors. Transcriptomic analysis confirmed that the erythroid GRN was modulated by pan-ErbB inhibitors, but not by Lapatinib, a dual EGFR/ErbB2 inhibitor (Supplementary Fig. 14a,b). Although cells were treated between the intermediate and late erythroblast stages, early network cluster genes (G1; Supplementary Fig. 14c) were significantly decreased, suggesting that ErbB signaling plays a role during multiple stages of differentiation. Analysis of pathways dysregulated by pan-ErbB inhibition revealed upregulation of P53 signaling (Fig. 5a), with concomitant downregulation of mitotic and proliferative pathways (Fig. 5b). Consistent with our prior computational analysis, these data connect ErbB signaling with P53 and proliferative pathways in human erythropoiesis.

As mechanistic analyses identified the Wnt pathway as a putative downstream target of ErbB4 (Fig. 5b), we exploited the pharmacologic accessibility of this pathway to enhance erythropoiesis *in vitro*. A critical barrier to blood generation as a cell-based biotechnology stems from a block in erythroid maturation from iPS cells (iRBCs), often requiring the use of feeder cells. We promoted maturation of iRBCs in a feeder-free system using bioprocess-compatible hematopoietic progenitors, which undergo continuous expansion under doxycycline-induced overexpression of 5 transcription factors (Fig. 5c)<sup>35</sup>. Activation of Wnt signaling via the agonist CHIR99021 increased the maturation of iRBCs, resulting in a 1.8 fold ( $p=4.8 \times 10^{-5}$ ) increase in GlyA<sup>+</sup>CD71<sup>-</sup> orthochromatic erythroblasts (Fig. 5d). Concomitantly, cells decreased in size with an increased nuclear-to-cytoplasmic ratio (Fig. 5e). Collectively, these data demonstrate that systems-level identification of druggable signaling pathways in developmental processes, such as erythropoiesis, is directly applicable to stem cell biomanufacturing and regenerative cell therapies.



## DISCUSSION

Here we establish the utility of systems-level analytics to elucidate biological processes that mediate dynamic stem cell and developmental transitions. Our computational pipeline provides a roadmap for the derivation of network models that connect sparse gene signatures with corresponding, yet disparate, biological processes, to capture the multi-factorial nature of cell state transitions. With cell engineering in hematopoiesis as an example, we highlight how to connect critical elements (e.g., LASSO gene signatures) to pathways/processes (e.g., networks derived via PCSF and correlation). Our network models suggested and we experimentally confirmed a previously unanticipated role for ErbB4 in hematopoiesis.

Our advanced pipeline integrates network topological architecture with pseudotemporal information to provide multiple layers of information about cell differentiation, which is complementary to purely trajectory-based algorithms<sup>7-10</sup> and highlights the changing roles of transcriptional regulators across dynamic stages of development. Moreover, the LASSO feature reduction as a foundation for network modeling ensures that the resulting models are informed by genes most vital to distinguishing divergent cell states. In contrast to traditional differential gene expression approaches<sup>19-21</sup>, LASSO produces a sparse, sharply focused gene set and when combined with PCSF produces a signaling network comprised of branches associated with distinct biological processes. Together, these approaches provide a more global depiction of the systems-level processes associated with cell fate transitions.

By applying our pipeline to study hematopoietic specification, we established a novel role for ErbB4 signaling in erythropoiesis in multiple *in vitro* and *in vivo* models. Many of our computational approaches did not directly identify ErbB4; however, network propagation from our maturation signature repeatedly identified ErbB ligands and ErbB-associated signaling, including MAPK/ERK, mitotic processes, P53, and apoptosis<sup>36,37</sup>. This highlights the need for future development of unsupervised metrics to prioritize candidates from aggregate data, which currently requires expert knowledge as an integral part of the process. Although there were no annotated processes enriched within the reticulocyte gene cluster, it included the NMDA receptor, GRIN3B, which is commonly implicated, along with ErbB4, in neurological development<sup>38</sup> and pathophysiology<sup>39</sup>. Interestingly, anemia is a common side effect of antipsychotic drugs<sup>40</sup> and studies of glutamate-mediated ion channels supports their functional role in erythropoiesis<sup>41</sup>. This opens the possibility of new avenues of crosstalk between neurological and hematopoietic systems, akin to the regulation of hematopoietic stem cell (HSC) production by the central nervous system<sup>42</sup>. Our dynamic analyses also revealed that oxidative stress pathways peak at the late erythroblast stage; ErbB4 is a known stress responsive pathway in the heart<sup>43</sup> and abrogates oxidative damage in the brain<sup>44</sup>. Although a recent meta-analysis of GWAS data identified neuregulin-4 (NRG4), an ErbB4-specific ligand, as a putative locus in aberrant human RBC phenotypes<sup>45</sup>, the pathway has not been previously characterized in erythropoiesis.

Cell engineering has broadly focused on inducing transcription factors as the emissaries of phenotype. To this end, CellNet successfully predicts candidate and aberrant transcription factors. However, even the most-studied form of reprogramming, induced pluripotency, remains exquisitely sensitive to culture conditions and relies on signaling molecules, such as

bFGF<sup>46</sup>. The advanced CellNet pipeline demonstrated here allows transcriptional targets to be complemented with druggable pathways. We demonstrate that the downstream ErbB signaling pathway can be exploited as a druggable target for more robust production of RBCs from an iPS-derived, bioprocess compatible, progenitor. Such multi-level reprogramming strategies may be especially beneficial for establishing and maintaining elusive populations, such as HSCs. Although engraftable HSCs can be generated with transcription factors alone<sup>47</sup>, reprogramming is enhanced by perturbation of developmental pathways, such as TGF $\beta$  and BMP4<sup>48</sup>. Similarly, AKT-activated endothelial cells support self-renewal and maintenance of HSCs through angiocrine factors<sup>49</sup>. The prevalence of growth factor supplementation and co-cultures across hematopoietic differentiation protocols further highlights the need to identify and recapitulate<sup>50</sup> cell-extrinsic signals.

## ONLINE METHODS

### GRN reconstruction & CellNet analytics

164 erythroid Affymetrix microarrays (Supplementary Tables 1&2) from the HGU133plus2 platform were acquired from the Gene Expression Omnibus (GEO) and compiled with the original human CellNet compendium<sup>51</sup>. Microarrays were preprocessed, the global gene regulatory network (GRN) was calculated via the Context Likelihood of Relatedness (CLR)<sup>52</sup> inference algorithm, and subnetworks were detected via InfoMap<sup>53</sup> community detection, as previously described<sup>51,54</sup>. Unless specified, all high dimensionality data analytics were accomplished using the R computational environment (version 3.2.2), with specified packages from CRAN and Bioconductor. All graphical representations and network analytics were visualized and calculated with the igraph package. Cell and tissue specific GRNs were established via enrichment using the chi-squared statistical test. As implemented in CellNet, a random forest classifier was trained based on the GRN for each cell type and trained with a randomly selected subset comprised of approximately 50% of the microarrays. The classification performance was then evaluated on the remaining independent subset of microarrays. The sensitivity (true positive divided by the sum of true positives and false negatives) at a false positive rate of 5% was calculated as a metric to evaluate classifier accuracy. The GRN score, defined as the weighted mean of expression Z-score, was calculated as previously described<sup>51</sup>.

### PCA and trajectory establishment

To evaluate GRN dynamics, the 235 genes within the erythroid GRN were mean centered and scaled to unit variance prior to dimensionality reduction via decomposition of the multivariate dataset into Principal Component (PC) space via the *prcomp* command from the stats package. GSEA<sup>55</sup> enrichment scores were calculated for each sample as a pre-ranked list, relative to the population mean, and correlations with respect to the PC axes were calculated via Pearson's *r*. All GSEA analyses were run using the Hallmark datasets in the Molecular Signatures Database (MSigDB)<sup>56</sup>. Unsupervised clustering of both microarrays and genes were calculated within the PC space via Gaussian Mixture Model-based methods from the mclust package.

For global analysis of gene dynamics, differential genes were first identified via the following criteria: 1) expression above a minimum threshold of 3.5, 2) variation across the dataset using an interquartile range greater than 0.75, and 3) significance between clusters via Bonferroni adjusted ANOVA. The remaining 1788 differentially expressed genes were clustered via K-means into K=12 groups, with the K determined by calculating Bayesian inference criteria (BIC). All enrichment analyses were conducted using standard chi-squared or Fisher statistics on Gene Ontology (GO)<sup>57,58</sup> defined biological processes.

### LASSO & network propagation

The Least Absolute Shrinkage and Selection Operator (LASSO)<sup>59</sup> was calculated within the *glmnet* package<sup>60</sup> using a binomial classification of microarrays from clusters C5 and C6 (Fig. 1b), based upon the differential genes (1788 genes) defined via criteria above. The value of lambda ( $\lambda=0.009$ ), calculated via *cv.glmnet*, was chosen by minimizing the mean square error (MSE), and corresponded to a signature of 27 genes. To validate the LASSO model, a 2-component PLSDA model (*mixOmics* package<sup>61</sup>) was built based upon the 27 genes and used to predict the binary classifications into clusters C5 and C6 (Supplementary Fig. 6). The calibration accuracy and error rates were calculated by the “Leave One Out” (LOO) method and compared to random models by: 1) shuffling the classifications and 2) selecting 27 random genes. In each case, the performance of random models was determined based on the average of 1000 permutations.

To expand the network without overfitting, three additional models were built based upon: 1) CellNet transcriptional regulatory networks, 2) protein-protein signaling networks, and 3) coexpression networks (Fig. 2b). The CellNet first order network was derived from the amalgamation of all first order connections to the 27-gene signature within the global CellNet GRN. P-values corresponding to the overrepresentation of each network regulator were calculated based upon Fisher’s test comparing the connections within the first order network to those in the global GRN and corrected for multiple hypothesis testing. Regulatory modules within the first order network were determined using the walktrap community algorithm.

The protein-protein interaction network was derived via querying the STRING database (version 9.0)<sup>62,63</sup> using the Prize Collecting Steiner Forest (PCSF) algorithm, as previously described<sup>64</sup>. Low confidence interactions with an edge score,  $s(e) < 0.5$  were removed and the cost was calculated as  $1 - s(e)$ . For increased robustness, noise was added to the edge cost and the resulting network was the amalgamation of 10 iterations, as previously described<sup>65</sup>. PCSF parameters, including  $\mu$  (node degree penalty),  $\omega$  (number of trees) and  $\beta$  (node prize scaling), were varied to demonstrate that the network generation was robust across a range of values (Supplementary Fig. 8). Larger networks (increasing  $\omega$  and  $\beta$ ), exhibited decreased density and centralization, with an increase in the number of significant ( $p < 0.05$ ) Gene Ontology annotations. This suggests that growing larger networks contributes toward the inclusion of discrete but cohesive biological processes, rather than adding random, unrelated genes/proteins. Moreover, smaller networks (Supplementary Fig. 8c–e) also exhibit common regulatory nodes consistent with our complete LASSO signaling network (Fig. 2f)

The coexpression network was determined by calculating the Pearson's  $r$  between each of the 27 signature genes and all other genes across clusters C2, C4, C5 and C6 (Fig. 1b). An absolute cutoff of 0.9 was selected for network reconstruction, based upon the "elbow" of network size over the range of thresholds (Supplementary Fig. 9a). The network parameters, as well as Gene Ontology annotations were calculated across the full range of cutoff thresholds and representative networks spanning the range exhibit similar features in terms of predicted transcription factors (ChEA/ENCODE), kinases (LINCS L1000) and ligands (LINCS L1000) (Supplementary Fig. 9b–e). Enrichment analyses for the all networks were queried via Gene Ontology (GO)<sup>57,58</sup> and Enrichr<sup>66</sup> in their native implementations.

### Statistical analyses

All statistical analyses were calculated in R, using two-sided, unpaired t-test, ANOVA, or Fisher's exact test. Data are presented as standard boxplots representing the median and ranging from the 25th to 75th percentiles, with the whiskers extending to 1.5\*IQR. The sample sizes represent a minimum of three independent replicates, corresponding to distinct experiments and/or parallel biological replicates (e.g. animals or cell cultures). The exact replicate numbers and statistical tests are specified in the figure legends.

### Human CD34+ RBC differentiation

Human CD34<sup>+</sup> progenitors derived from mobilized peripheral blood (AllCells) were expanded for 4 days in StemSpan SFEM (STEMCELL Technologies) with the addition of IL3 (10 ng/mL), IL6 (50 ng/mL), TPO (50 ng/mL), SCF (50 ng/mL) and Flt3 (50 ng/mL). Unless specified, all cytokines were from PeproTech. Erythroid differentiation was accomplished using a previously published, 3-stage protocol<sup>67</sup>. Briefly, all stages of differentiation consisted of a basal erythroid differentiation medium (EDM) comprised of: IMDM with 15% FBS, 1% BSA, 2 mM L-glutamine, 500  $\mu$ g/mL holo-transferrin, and 10  $\mu$ g/mL insulin. Stage 1 consists of EDM plus the addition of dexamethasone (1  $\mu$ M),  $\beta$ -estradiol (1  $\mu$ M), IL3 (5 ng/mL), SCF (100 ng/mL) and EPO (6U) for 5 days (days 0–5). Stage 2 consists of EDM plus the addition of SCF (50 ng/mL) and EPO (6U) for 4 days (days 5–9). Stage 3 consists of EDM plus the addition of EPO alone (2U) for 8 days (days 9–17). At all stages, cells are cultured in 24 well plates in 1 mL of media. Cell number seeded at the beginning of stages 1, 2, and 3 are:  $10^5$ ,  $2 \times 10^5$ , and  $5 \times 10^5$ /well.

### iPS-5F generation and RBC differentiation

Human iPS-5F cells were generated as previously described<sup>68,69</sup> from MSC-iPS<sup>70</sup> obtained from the Boston Children's Hospital Human Embryonic Stem Cell Core (hESC) and verified by immunohistochemistry for pluripotency markers, teratoma formation and karyotyping. Briefly, iPS cells were differentiated as embryoid bodies using a hematopoietic induction protocol<sup>71</sup> and CD34<sup>+</sup> cells were sorted from bulk embryoid body culture by magnetic activated cell sorting (MACS) using human CD34 microbeads (Miltenyi Biotec), as per the manufacturer's instructions. The embryoid body progenitors were seeded on retronectin-coated ( $10 \mu$ g  $\text{cm}^{-2}$ ) 96-well plates ( $2 \times 10^4$ – $5 \times 10^4$  cells per well) in SFEM (StemCell Technologies) containing 50 ng  $\text{ml}^{-1}$  SCF, 50 ng  $\text{ml}^{-1}$  FLT3, 50 ng  $\text{ml}^{-1}$  TPO (all R&D Systems), 50 ng  $\text{ml}^{-1}$  IL-6 and 10 ng  $\text{ml}^{-1}$  IL-3 (both from Peprotech) and infected with 5F lentiviral particles. Lentiviral particles for the 5F plasmids (HOXA9, ERG, RORA, SOX4

and MYB cloned into pInducer-21 doxycycline-inducible vector) were produced by transfecting 293T-17 cells (ATCC) with third-generation packaging plasmids. The multiplicity of infection (MOI) for each factor was: ERG MOI = 5, HOXA9 MOI = 5, RORA MOI = 3, SOX4 MOI = 3, MYB MOI = 3. Following 24 hours of infection, 5F cells were cultured in SFEM with 50 ng ml<sup>-1</sup> SCF, 50 ng ml<sup>-1</sup> FLT3, 50 ng ml<sup>-1</sup> TPO, 50 ng ml<sup>-1</sup> (all R&D Systems) IL-6, and 10 ng ml<sup>-1</sup> IL-3 (PeproTech) and 2 µg ml<sup>-1</sup> doxycycline (Dox; Sigma). Cultures were maintained at a density of <math>1 \times 10^6</math> cells ml<sup>-1</sup>, and the medium was changed every 3–4 days.

RBC differentiation from iPS-5F followed a slightly modified protocol that was previously optimized for translational approaches aimed at transfusion of in vitro-derived RBCs.<sup>72</sup> In this protocol, the EDM was instead comprised of: IMDM with 5% inactivated plasma (Solvent Detergent Pooled Plasma AB from the Rhode Island Blood Center), 2 mM L-glutamine, 330 µg/mL holo-transferrin, and 10 µg/mL insulin, 2 IU/mL heparin (Sigma) and 3 IU/mL EPO. Stage I (days 0–7) was plated at 1–3 × 10<sup>5</sup> cells/mL and supplemented with 10 µM hydrocortisone, 100 ng/mL SCF and 5 ng/mL IL-3. Stage II (days 7–11) was plated at 1–3 × 10<sup>5</sup> cells/mL and supplemented with 100 ng/mL SCF. Stage III (days 11–18) was plated at 1 × 10<sup>6</sup> cells/mL in the basal EDM. All analyses were conducted at day 18 of differentiation and CHIR99021 (3 µM) was added throughout stage III (days 9 and 13).

### Flow cytometry & cell sorting

Human erythropoiesis, including differentiation from BM CD34+ cells and native bone marrow samples, was analyzed with the following antibody panel: CD71 PE (M-A712; BD), and CD235a/Glycophorin A PE- Cy7 (11E4B-7–6; Coulter) or CD235a/Glycophorin A FITC (11E4B-7–6; Coulter). Mouse erythropoiesis from Neratinib treated and HER4<sup>heart</sup> mice was analyzed with the following antibody panel: mCD71 FITC (C2; BD), mTer119 PE-Cy5 (Ter-119; eBioscience). All staining was performed with <math>1 \times 10^6</math> cells per 100 µL staining buffer (PBS + 2% FBS) with 1:100 dilution of each antibody for 30 min at RT in dark. Compensation was performed by automated compensation with anti-mouse Igk and negative beads (BD). Acquisition was performed on a BD Fortessa cytometer and all sorting was performed on a BD FACS Aria II cell sorter using a 70-mm nozzle. Gating strategies are depicted in Supplementary Fig. 15.

### Inhibitors

All inhibitors were added to cell cultures at 1 µM on days 9 and 13 of differentiation, corresponding to the beginning and middle of stage 3 (supplemented with EPO only, as described above). DMSO was used for a vehicle control in all cell culture studies. Details on ordering information and affinities are provided in Supplementary Table 13.

### RNA-sequencing

RNA was extracted after 24 hours of incubation with ErbB inhibitors (day 10 of erythroid differentiation) using Trizol reagent (Invitrogen) and the RNeasy Plus kit (Qiagen). Quality of RNA was monitored via QC for high RIN values and low levels of DNA contamination. RNA-seq libraries were prepared using the SMARTseq v4 kit as per manufacturer's protocol with 10 ng input RNA. Libraries were sequenced using the 200 cycle paired-end kit on the

Illumina HiSeq2500 system. RNA-seq reads were analyzed with the Tuxedo Tools following a standard protocol on the Harvard Medical School Orchestra Cluster. Reads were mapped with TopHat version 2.1.0 and Bowtie2 version 2.2.4 with default parameters against build hg19 of the human genome, and build hg19 of the RefSeq human genome annotation. Samples were quantified with the Cufflinks package version 2.2.1. Differential expression was performed using Cuffdiff with default parameters.

## PCR

RNA was extracted as described above and cDNA was synthesized using the SuperScript™ VILO™ cDNA Synthesis Kit (Thermo), per manufacturer's instructions. Real time PCR was run using SYBR green technology with QuantiTect primers for the ErbB receptor family (Qiagen) on the QuantStudio Flex Real-Time PCR System.

## Zebrafish studies

Zebrafish were maintained according to institutional animal care and use committee-approved protocols. The *Tg(globin:eGFP)* line was provided by L. I. Zon, Children's Hospital, Harvard Medical School, Boston, MA. MOs (GeneTools) were microinjected at the one-cell stage as described previously<sup>73</sup>. ErbB4 MOs were generated from previously published sequences<sup>74</sup>. Embryos were harvested at 48–56 hours post fertilization (hpf) and were processed with matched sibling controls for o-dianisidine staining and evaluation of globin:eGFP intensity. Staining intensity was categorized as low, medium or high within that experiment, as previously described<sup>75</sup>; effects were independently confirmed by other lab members.

## Mouse studies

All mice were housed in pathogen-free animal facilities, and all experiments were performed with the approval of the Animal Care and Use Committee at Harvard Medical School and Dana-Farber Cancer Institute and/or the BCH animal care committee. At least n=3 animals were used per cohort, based on previous studies. For drug treatments, mice were assigned randomly to groups and not blinded. Neratinib was delivered to B6 albino mice via oral gavage at 60 mg/kg daily for 1 week. Hydroxypropyl methylcellulose (HPMC) was used as a vehicle control for Neratinib in mouse treatments.

## Reporting Summary

Further information on research design is available in the Nature Research Reporting Summary linked to this article.

## Supplementary Material

Refer to Web version on PubMed Central for supplementary material.

## Acknowledgements.

The authors thank Dr. Gabriel Corfas (University of Michigan) for sharing the ErbB4<sup>-/-</sup>HER4heart mutant mice, which were generated in 2003 by Martin Gassmann (University of Basel) and colleagues<sup>34</sup> and Dr. Leonard I. Zon (Boston Children's Hospital) for the *globin:eGFP* transgenic fish. The authors also thank Pinar Eser for the ErbB

inhibitor library and Tolulope Rosanwo for cells and reagents, as well as Ronald Mathieu and the BCH Flow Cytometry Core, Jihan Osborne, Brian Joughin, Jishnu Das, and Annelien Zweemer for technical advice. FUNDING: This work is supported by grants from the NIH NIGMS (R01-GM081336), NIH NIDDK (R24-DK092760, R24-DK49216) and NHLBI Progenitor Cell Translation Consortium (U01HL134812); NHLBI R01-HL04880, NIH R24-OD017870-01 and NIGMS R01-GM069668. M.A.K. is supported by a NIH T32 Training Grant from BWH Hematology. L.T.V. was supported by the NSF Graduate Research Fellowship. J.M.F. is supported by a NIH T32 Training Grant from the NHLBI. T.E.N. is a Leukemia and Lymphoma Society Scholar. G.Q.D. is an associate member of the Broad Institute and was supported by the Howard Hughes Medical Institute and the Manton Center for Orphan Disease Research.

## REFERENCES

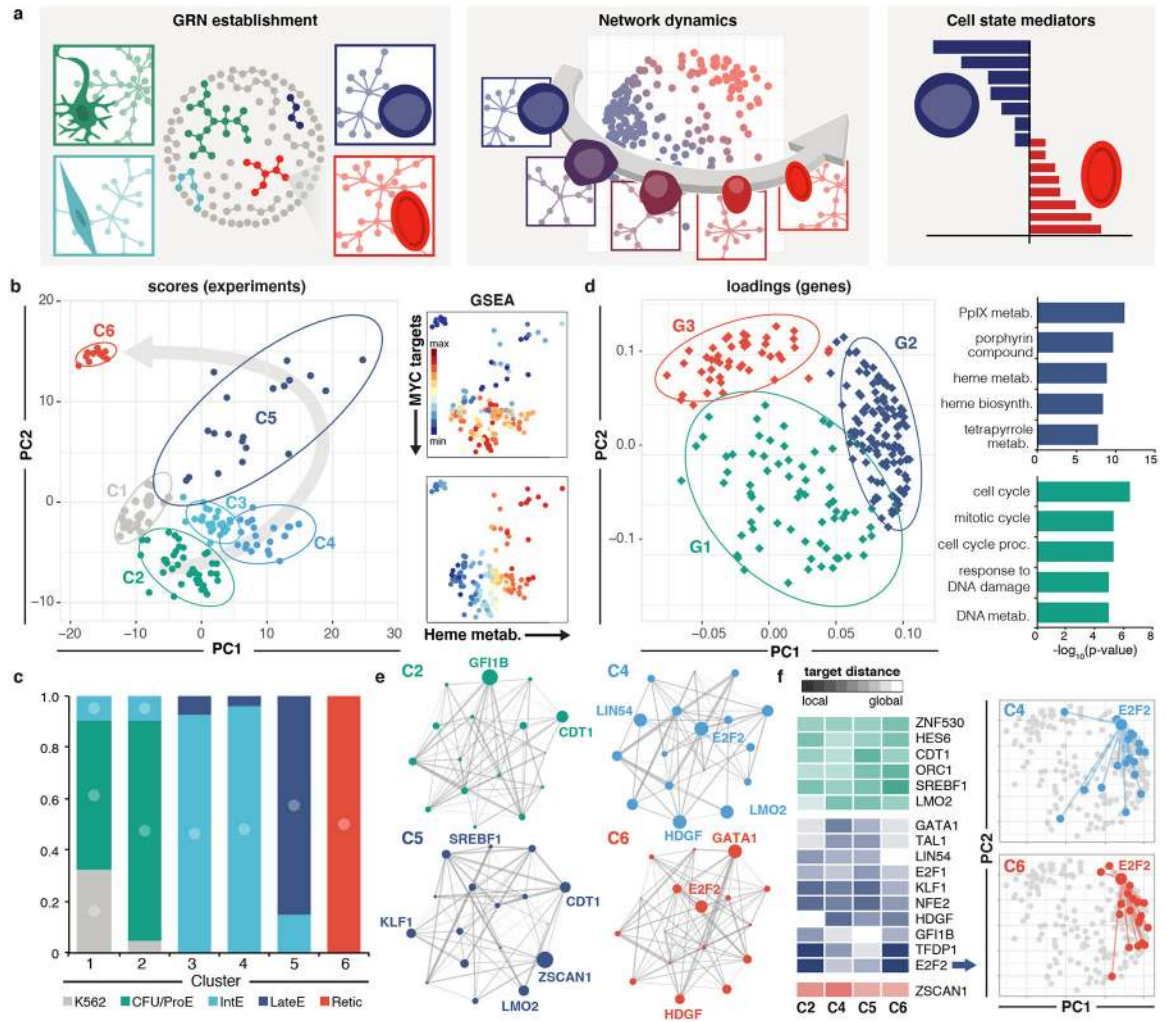
1. Westerhoff HV & Palsson BO The evolution of molecular biology into systems biology. *Nat Biotechnol* 22, 1249–1252 (2004). [PubMed: 15470464]
2. Morris SA & Daley GQ A blueprint for engineering cell fate: current technologies to reprogram cell identity. *Cell Res* 23, 33–48 (2013). [PubMed: 23277278]
3. Hanna JH, Saha K & Jaenisch R Pluripotency and cellular reprogramming: facts, hypotheses, unresolved issues. *Cell* 143, 508–525 (2010). [PubMed: 21074044]
4. Cahan P et al. CellNet: Network Biology Applied to Stem Cell Engineering. *Cell* 158, 903–915 (2014). [PubMed: 25126793]
5. Morris SA et al. Dissecting Engineered Cell Types and Enhancing Cell Fate Conversion via CellNet. *Cell* 158, 889–902 (2014). [PubMed: 25126792]
6. Radley AH et al. Assessment of engineered cells using CellNet and RNA-seq. *Nat Protoc* 12, 1089–1102 (2017). [PubMed: 28448485]
7. Setty M et al. Wishbone identifies bifurcating developmental trajectories from single-cell data. *Nat Biotechnol* 34, 637–645 (2016). [PubMed: 27136076]
8. Shin J et al. Single-Cell RNA-Seq with Waterfall Reveals Molecular Cascades underlying Adult Neurogenesis. *Cell Stem Cell* 17, 360–372 (2015). [PubMed: 26299571]
9. Trapnell C et al. the dynamics and regulators of cell fate decisions are revealed by pseudotemporal ordering of single cells. *Nat Biotechnol* 1–11 (2014). doi:10.1038/nbt.2859 [PubMed: 24406907]
10. Lummertz da Rocha E et al. Reconstruction of complex single-cell trajectories using CellRouter. *Nat Commun* 9, 892 (2018). [PubMed: 29497036]
11. Dzierzak E & Philipsen S Erythropoiesis: development and differentiation. *Cold Spring Harb Perspect Med* 3, a011601–a011601 (2013). [PubMed: 23545573]
12. Tsai FY & Orkin SH Transcription factor GATA-2 is required for proliferation/survival of early hematopoietic cells and mast cell formation, but not for erythroid and myeloid terminal differentiation. *Blood* 89, 3636–3643 (1997). [PubMed: 9160668]
13. Cantor AB & Orkin SH Transcriptional regulation of erythropoiesis: an affair involving multiple partners. *Oncogene* 21, 3368–3376 (2002). [PubMed: 12032775]
14. da Cunha AF et al. Global gene expression reveals a set of new genes involved in the modification of cells during erythroid differentiation. *Cell Prolif* 43, 297–309 (2010). [PubMed: 20546246]
15. Ding K et al. Genetic Loci implicated in erythroid differentiation and cell cycle regulation are associated with red blood cell traits. *Mayo Clin. Proc* 87, 461–474 (2012). [PubMed: 22560525]
16. Li J et al. Isolation and transcriptome analyses of human erythroid progenitors: BFU-E and CFU-E. *Blood* 124, 3636–3645 (2014). [PubMed: 25339359]
17. Rylski M et al. GATA-1-mediated proliferation arrest during erythroid maturation. *Mol Cell Biol* 23, 5031–5042 (2003). [PubMed: 12832487]
18. Goh S-H et al. The human reticulocyte transcriptome. *Physiol. Genomics* 30, 172–178 (2007). [PubMed: 17405831]
19. Ideker T, Ozier O, Schwikowski B & Siegel AF Discovering regulatory and signalling circuits in molecular interaction networks. *Bioinformatics* 18 Suppl 1, S233–40 (2002). [PubMed: 12169552]
20. Langfelder P, Mischel PS & Horvath S When is hub gene selection better than standard meta-analysis? *PLoS ONE* 8, e61505 (2013). [PubMed: 23613865]
21. Barabási A-L, Gulbahce N & Loscalzo J Network medicine: a network-based approach to human disease. *Nat Rev Genet* 12, 56–68 (2011). [PubMed: 21164525]

22. Takahashi K et al. Induction of pluripotent stem cells from adult human fibroblasts by defined factors. *Cell* 131, 861–872 (2007). [PubMed: 18035408]
23. Yu J et al. Induced pluripotent stem cell lines derived from human somatic cells. *Science* 318, 1917–1920 (2007). [PubMed: 18029452]
24. Tuncbag N et al. Simultaneous reconstruction of multiple signaling pathways via the prize-collecting steiner forest problem. *J. Comput. Biol* 20, 124–136 (2013). [PubMed: 23383998]
25. Robert-Moreno A, Espinosa L, Sanchez MJ, la Pompa, de JL & Bigas A The notch pathway positively regulates programmed cell death during erythroid differentiation. *Leukemia* 21, 1496–1503 (2007). [PubMed: 17476283]
26. Watanabe S et al. Loss of a Rho-regulated actin nucleator, mDia2, impairs cytokinesis during mouse fetal erythropoiesis. *Cell Rep* 5, 926–932 (2013). [PubMed: 24239357]
27. Tanno T et al. High levels of GDF15 in thalassemia suppress expression of the iron regulatory protein hepcidin. *Nat Med* 13, 1096–1101 (2007). [PubMed: 17721544]
28. Jacquelin A et al. Apoptosis and erythroid differentiation triggered by Bcr-Abl inhibitors in CML cell lines are fully distinguishable processes that exhibit different sensitivity to caspase inhibition. *Oncogene* 26, 2445–2458 (2007). [PubMed: 17043649]
29. Zhou BP et al. HER-2/neu induces p53 ubiquitination via Akt-mediated MDM2 phosphorylation. *Nat Cell Biol* 3, 973–982 (2001). [PubMed: 11715018]
30. Le XF et al. Heregulin-induced apoptosis is mediated by down-regulation of Bcl-2 and activation of caspase-7 and is potentiated by impairment of protein kinase C alpha activity. *Oncogene* 20, 8258–8269 (2001). [PubMed: 11781840]
31. Holbro T et al. The ErbB2/ErbB3 heterodimer functions as an oncogenic unit: ErbB2 requires ErbB3 to drive breast tumor cell proliferation. *Proc Natl Acad Sci USA* 100, 8933–8938 (2003). [PubMed: 12853564]
32. Lee H-Y et al. PPAR- $\alpha$  and glucocorticoid receptor synergize to promote erythroid progenitor self-renewal. *Nature* 522, 474–477 (2015). [PubMed: 25970251]
33. Orkin SH & Zon LI Hematopoiesis: an evolving paradigm for stem cell biology. *Cell* 132, 631–644 (2008). [PubMed: 18295580]
34. Tidcombe H et al. Neural and mammary gland defects in ErbB4 knockout mice genetically rescued from embryonic lethality. *Proc Natl Acad Sci USA* 100, 8281–8286 (2003). [PubMed: 12824469]
35. Doulatov S et al. Induction of Multipotential Hematopoietic Progenitors from Human Pluripotent Stem Cells via Respecification of Lineage-Restricted Precursors. *Cell Stem Cell* 13, 459–470 (2013). [PubMed: 24094326]
36. Naresh A et al. The ERBB4/HER4 intracellular domain 4ICD is a BH3-only protein promoting apoptosis of breast cancer cells. *Cancer Res* 66, 6412–6420 (2006). [PubMed: 16778220]
37. Bersell K, Arab S, Haring B & Kühn B Neuregulin1/ErbB4 signaling induces cardiomyocyte proliferation and repair of heart injury. *Cell* 138, 257–270 (2009). [PubMed: 19632177]
38. Li B, Woo R-S, Mei L & Malinow R The neuregulin-1 receptor erbB4 controls glutamatergic synapse maturation and plasticity. *Neuron* 54, 583–597 (2007). [PubMed: 17521571]
39. Hahn C-G et al. Altered neuregulin 1-erbB4 signaling contributes to NMDA receptor hypofunction in schizophrenia. *Nat Med* 12, 824–828 (2006). [PubMed: 16767099]
40. Flanagan RJ & Dunk L Haematological toxicity of drugs used in psychiatry. *Hum Psychopharmacol* 23 Suppl 1, 27–41 (2008). [PubMed: 18098216]
41. Hänggi P et al. Functional plasticity of the N-methyl-d-aspartate receptor in differentiating human erythroid precursor cells. *Am J Physiol, Cell Physiol* 308, C993–C1007 (2015). [PubMed: 25788577]
42. Kwan W et al. The Central Nervous System Regulates Embryonic HSPC Production via Stress-Responsive Glucocorticoid Receptor Signaling. *Cell Stem Cell* 19, 370–382 (2016). [PubMed: 27424782]
43. Kuramochi Y et al. Cardiac endothelial cells regulate reactive oxygen species-induced cardiomyocyte apoptosis through neuregulin-1beta/erbB4 signaling. *J Biol Chem* 279, 51141–51147 (2004). [PubMed: 15385548]



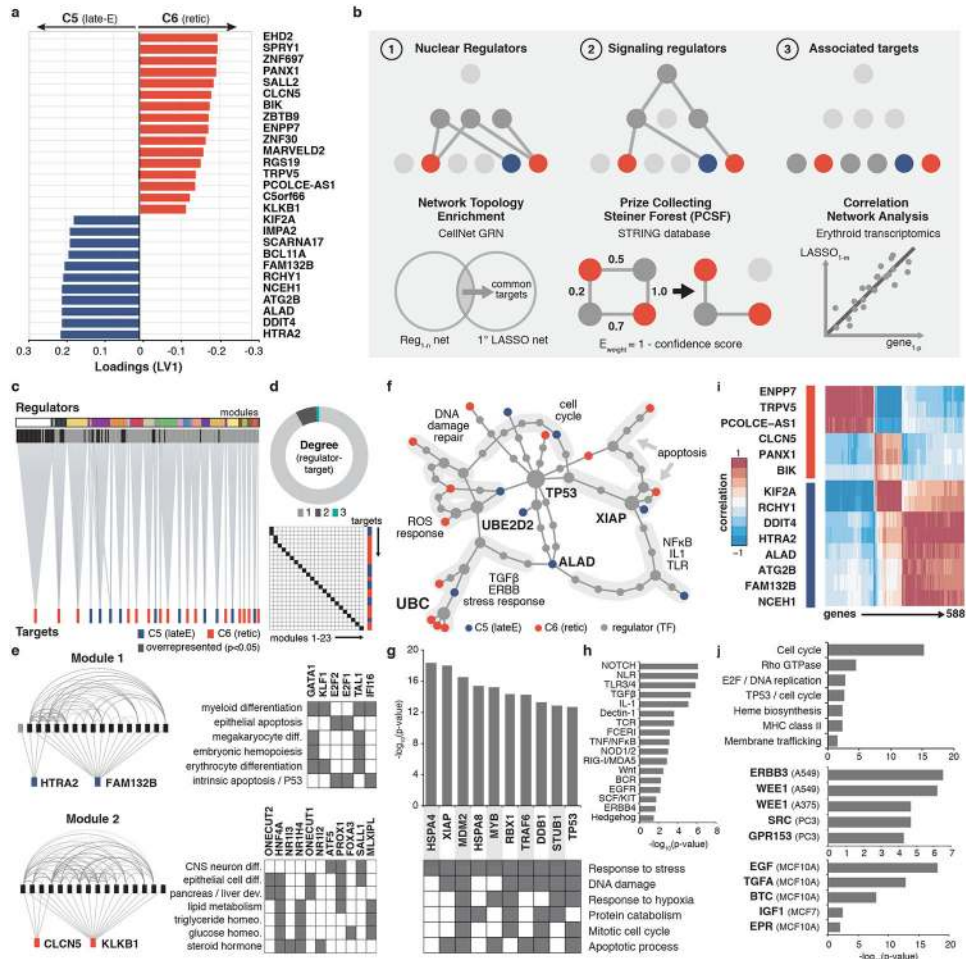
44. Xu Z et al. Neuroprotection by neuregulin-1 following focal stroke is associated with the attenuation of ischemia-induced pro-inflammatory and stress gene expression. *Neurobiol Dis* 19, 461–470 (2005). [PubMed: 16023588]
45. van der Harst P et al. Seventy-five genetic loci influencing the human red blood cell. *Nature* 492, 369–375 (2012). [PubMed: 23222517]
46. Ludwig TE et al. Derivation of human embryonic stem cells in defined conditions. *Nat Biotechnol* 24, 185–187 (2006). [PubMed: 16388305]
47. Sugimura R et al. Haematopoietic stem and progenitor cells from human pluripotent stem cells. *Nature* 545, 432–438 (2017). [PubMed: 28514439]
48. Lis R et al. Conversion of adult endothelium to immunocompetent haematopoietic stem cells. *Nature* 545, 439–445 (2017). [PubMed: 28514438]
49. Kobayashi H et al. Angiocrine factors from Akt-activated endothelial cells balance self-renewal and differentiation of haematopoietic stem cells. *Nat Cell Biol* 12, 1046–1056 (2010). [PubMed: 20972423]
50. Shukla S et al. Progenitor T-cell differentiation from hematopoietic stem cells using Delta-like-4 and VCAM-1. *Nat Methods* 14, 531–538 (2017). [PubMed: 28394335]
51. Cahan P et al. CellNet: Network Biology Applied to Stem Cell Engineering. *Cell* 158, 903–915 (2014). [PubMed: 25126793]
52. Faith JJ et al. Large-scale mapping and validation of *Escherichia coli* transcriptional regulation from a compendium of expression profiles. *PLoS Biol* 5, e8 (2007). [PubMed: 17214507]
53. Rosvall M & Bergstrom CT An information-theoretic framework for resolving community structure in complex networks. *Proc Natl Acad Sci USA* 104, 7327–7331 (2007). [PubMed: 17452639]
54. Radley AH et al. Assessment of engineered cells using CellNet and RNA-seq. *Nat Protoc* 12, 1089–1102 (2017). [PubMed: 28448485]
55. Subramanian A et al. Gene set enrichment analysis: a knowledge-based approach for interpreting genome-wide expression profiles. *Proc Natl Acad Sci USA* 102, 15545–15550 (2005). [PubMed: 16199517]
56. Liberzon A et al. The Molecular Signatures Database Hallmark Gene Set Collection. *Cell Syst* 1, 417–425 (2015). [PubMed: 26771021]
57. Ashburner M et al. Gene ontology: tool for the unification of biology. The Gene Ontology Consortium. *Nat Genet* 25, 25–29 (2000). [PubMed: 10802651]
58. Gene Ontology Consortium. Gene Ontology Consortium: going forward. *Nucleic Acids Res* 43, D1049–56 (2015). [PubMed: 25428369]
59. Tibshirani R Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society Series B (...)* (1996).
60. Friedman J, Hastie T & Tibshirani R Regularization Paths for Generalized Linear Models via Coordinate Descent. *J Stat Softw* 33, 1–22 (2010). [PubMed: 20808728]
61. Lê Cao K-A, Boitard S & Besse P Sparse PLS discriminant analysis: biologically relevant feature selection and graphical displays for multiclass problems. *BMC Bioinformatics* 12, 253 (2011). [PubMed: 21693065]
62. Snel B, Lehmann G, Bork P & Huynen MA STRING: a web-server to retrieve and display the repeatedly occurring neighbourhood of a gene. *Nucleic Acids Res* 28, 3442–3444 (2000). [PubMed: 10982861]
63. Szklarczyk D et al. The STRING database in 2017: quality-controlled protein-protein association networks, made broadly accessible. *Nucleic Acids Res* 45, D362–D368 (2017). [PubMed: 27924014]
64. Tuncbag N et al. Simultaneous reconstruction of multiple signaling pathways via the prize-collecting steiner forest problem. *J. Comput. Biol* 20, 124–136 (2013). [PubMed: 23383998]
65. Akhmedov M et al. PCSF: An R-package for network-based interpretation of high-throughput data. *PLoS Comput Biol* 13, e1005694 (2017). [PubMed: 28759592]
66. Chen EY et al. Enrichr: interactive and collaborative HTML5 gene list enrichment analysis tool. *BMC Bioinformatics* 14, 128 (2013). [PubMed: 23586463]

67. Lee H-Y et al. PPAR- $\alpha$  and glucocorticoid receptor synergize to promote erythroid progenitor self-renewal. *Nature* 522, 474–477 (2015). [PubMed: 25970251]
68. Doulatov S et al. Induction of Multipotential Hematopoietic Progenitors from Human Pluripotent Stem Cells via Respecification of Lineage-Restricted Precursors. *Cell Stem Cell* 13, 459–470 (2013). [PubMed: 24094326]
69. Vo LT et al. Regulation of embryonic haematopoietic multipotency by EZH1. *Nature* 553, 506–510 (2018). [PubMed: 29342143]
70. Park I-H et al. Reprogramming of human somatic cells to pluripotency with defined factors. *Nature* 451, 141–146 (2008). [PubMed: 18157115]
71. Chadwick K et al. Cytokines and BMP-4 promote hematopoietic differentiation of human embryonic stem cells. *Blood* 102, 906–915 (2003). [PubMed: 12702499]
72. Giarratana M-C et al. Proof of principle for transfusion of in vitro-generated red blood cells. *Blood* 118, 5071–5079 (2011). [PubMed: 21885599]
73. North TE et al. Prostaglandin E2 regulates vertebrate haematopoietic stem cell homeostasis. *Nature* 447, 1007–1011 (2007). [PubMed: 17581586]
74. Sato T et al. Neuregulin 1 Type II-ErbB Signaling Promotes Cell Divisions Generating Neurons from Neural Progenitor Cells in the Developing Zebrafish Brain. *PLoS ONE* 10, e0127360 (2015). [PubMed: 26001123]
75. Doulatov S et al. Drug discovery for Diamond-Blackfan anemia using reprogrammed hematopoietic progenitors. *Sci Transl Med* 9, (2017).



**Figure 1. GRN dynamics capture cell fate specification.**

(a) Schematic overview of computational pipeline. (b) Principal Component Analysis (PCA) scores of erythroid microarray datasets, with Gaussian Mixture Model (GMM) derived clusters. (c) Distribution of blinded manual classifications (K562 erythroleukemia cell line, early erythroblasts: CFU/ProE, intermediate erythroblasts: IntE, late/orthochromatic erythroblasts: LateE, and reticulocytes) across unsupervised computational clusters C1–C6 from part (b), with circular annotations indicative of  $p < 0.05$  by Fisher's exact test. (d) PCA loadings plot for the 235 genes in the erythroid GRN and Gene Ontology biological process enrichment for GMM gene clusters, with G1–G3 correlating with early erythroblasts (green), intermediate/late erythroblasts (blue), reticulocytes (red), respectively. (e) Visualization of regulators from dynamic networks calculated using CLR inference across each of the clusters from C2, C4, C5 and C6. Node size correlates with the degree (number of targets) and line width corresponds to the CLR Z-score (confidence of interaction) between regulators. (f) The average Euclidean distance in the PC loadings space was calculated for all targets of each regulator across the dynamic network topologies depicted in part (e), with representative plots demonstrating the target distribution for E2F2 across networks from erythroblasts (C4) and reticulocytes (C6).



**Figure 2. ErbB signaling is implicated in erythroid differentiation.**

(a) Gene signature distinguishing clusters C5 (late-E; blue) and C6 (reticulocyte; red) from Fig. 2a, as determined by Least Absolute Shrinkage and Selection Operator (LASSO). (b) schematic of network propagation approaches across different regulatory layers, including (1) nuclear regulators, (2), signaling regulators, and (3) associated targets, corresponding to networks from parts (c), (f) and (i), respectively. (c) First order connection network derived from connecting the LASSO signature genes within the CellNet global GRN. Significantly enriched regulators ( $p < 0.05$  by Fisher's test for connections in the network compared to the global GRN) are shown in black, with community-derived (i.e. locally high density) modules depicted above in varying colors and the LASSO genes in red and blue, corresponding to the representation in part (a). (d) The degree of all regulators and LASSO gene module membership. (e) Graphical representation of Modules 1 and 2, with significantly enriched regulators (black) and the associated Gene Ontology biological process annotations. (f) Protein-protein interaction (PPI) network built from gene signature using the STRING database and the Prize Collecting Steiner Tree (PCSF) algorithm. Non-LASSO nodes (Steiner nodes) are depicted with the size proportional to the degree and LASSO nodes (terminals) are depicted in red and blue, corresponding to the representation in part (a). Modules within the network are highlighted and annotated with significantly enriched biological processes. (g) P-values ranking node enrichment (Fisher's test for

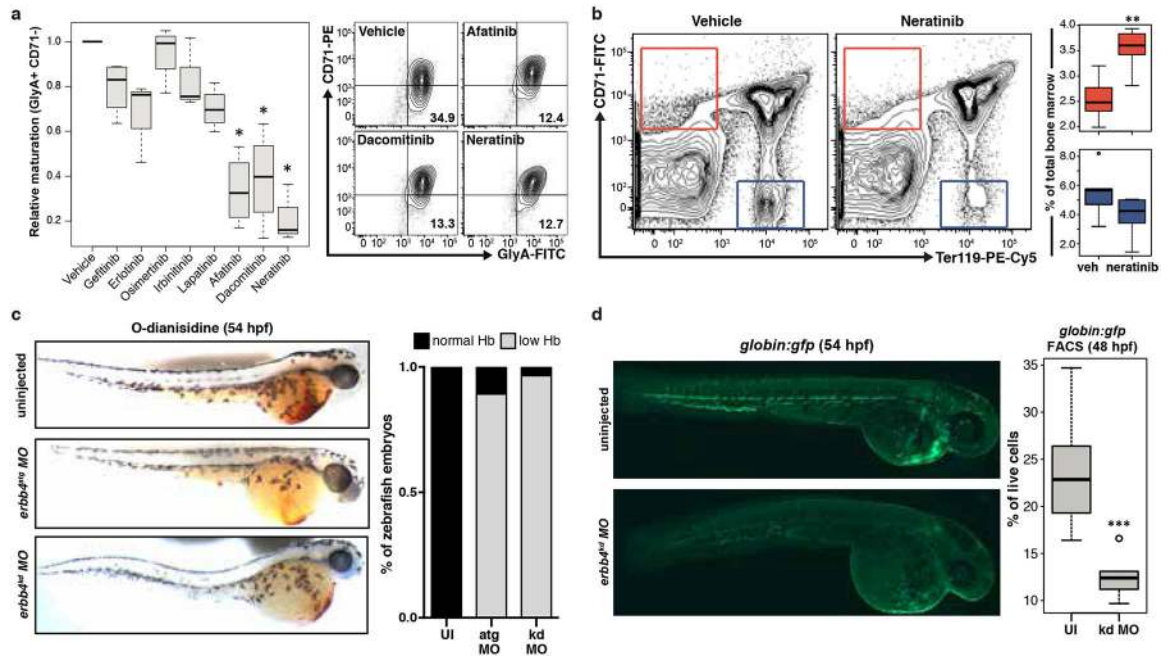
connections in the LASSO network relative to the full STRING network) and corresponding Gene Ontology annotations. (h) Enriched signaling pathways from the Reactome database. (i) Coexpression network comprising genes highly correlated  $r > |0.90|$  with the LASSO signature. (j) Enrichment analyses for Reactome processes, kinase perturbation (downregulated genes upon kinase knockdown from LINCS L1000) and ligand regulation (upregulated upon ligand stimulation from LINCS L1000).

Author Manuscript

Author Manuscript

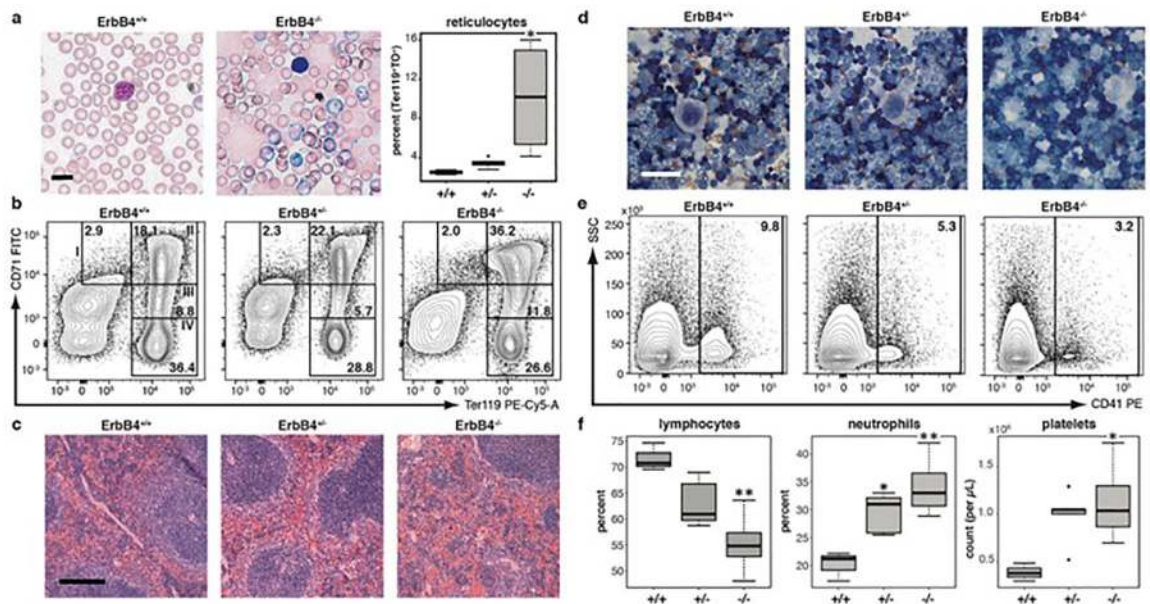
Author Manuscript

Author Manuscript



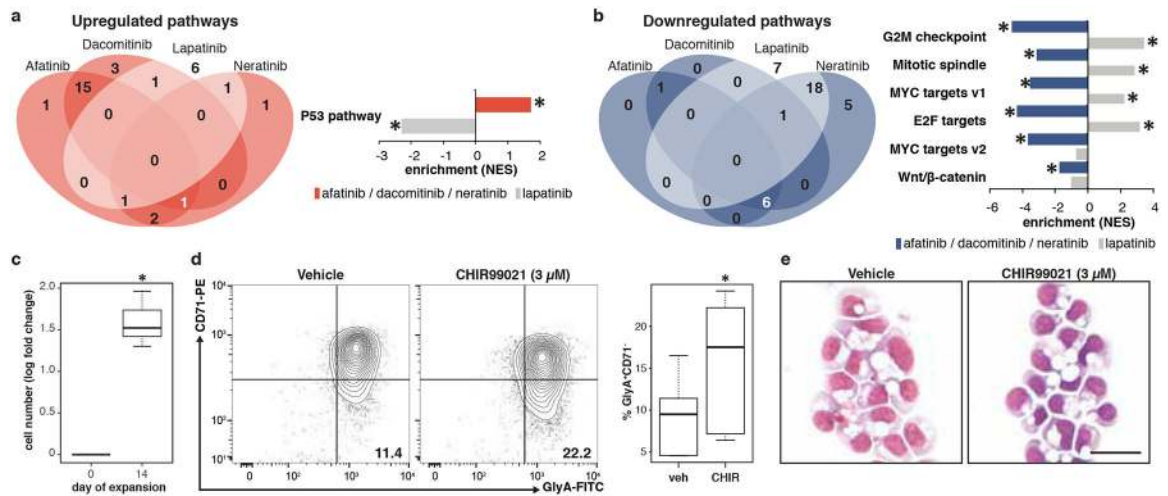
**FIGURE 3. ErbB4 is required for robust erythroid development.**

(a) Chemical screen with ErbB inhibitors during *in vitro* erythroid differentiation from bone marrow (BM) CD34<sup>+</sup> cells and quantified based upon the relative fractions of GlyA<sup>+</sup>CD71<sup>-</sup> cells, with representative flow cytometry plots. \* =  $p < 0.05$  compared to vehicle (DMSO), via ANOVA and post hoc Tukey HSD tests.  $n = 3-4$  independent experiments. (b) Early (CD71<sup>+</sup>GlyA<sup>-</sup>) and late erythroid (CD71<sup>-</sup>GlyA<sup>+</sup>) progenitor profiles in the bone marrow of mice ( $n = 5$  per group) treated with vehicle (0.5% hydroxypropyl methylcellulose; HPMC) or Neratinib (60 mg/kg) for 7 days. \*\* =  $p < 0.01$  by two-sided, unpaired t-test. (c) o-Dianisidine (3,3'-dimethoxybenzidine) staining of zebrafish 54 hpf after ErbB4 MO (Sato et al., PLoS One, 2015) injection (left). Quantification of the percentage of zebrafish embryos ( $n > 6$  across 3 clutches) with low hemoglobin (low o-dianisidine staining) relative to uninjected controls (right). (d) Transgenic globin-GFP zebrafish embryos ( $n > 6$  across 2 clutches) were visualized and analyzed via flow cytometry at 54 hpf and 48 hpf, respectively, after MO injection. \*\*\* =  $p < 0.001$  by two-sided, unpaired t-test. Scale bars correspond to 200  $\mu\text{m}$ .



**FIGURE 4. ErbB4 genetic deficiency leads to blood defects in the ErbB4<sup>-/-</sup>HER4<sup>heart</sup> mouse model.**

Morphology of the (a) peripheral blood from smears stained with May Grunwald Giemsa indicating the increased reticulocyte fraction (purple), quantified as Ter119<sup>+</sup> Thiazole Orange (TO)<sup>+</sup> cells (ErbB4<sup>+/+</sup> n=3; ErbB4<sup>-/-</sup> n=5; ErbB4<sup>+/-</sup> n=4; \* = p<0.05 compared to wild type by one way ANOVA). (b) Erythroid fractions within the spleen of wild type, heterozygous, and homozygous mice, demonstrating increased extramedullary erythropoiesis (gate II, CD71+Ter119+). (c) Morphology of the spleen stained with hematoxylin and eosin (H&E) showing dense red pulp. (d) Bone marrow composition from cytopspins stained with May Grunwald Giemsa demonstrating the presence of immature fraction hematopoietic fractions across multiple lineages. (e) Proportion of CD41<sup>+</sup> megakaryocytes in the spleen. (f) Quantification of lymphocytes, neutrophils, and platelets in the peripheral blood via Hemavet (ErbB4<sup>+/+</sup> n=3; ErbB4<sup>-/-</sup> n=5; ErbB4<sup>+/-</sup> n=4; \* = p<0.05, \*\* = p<0.01 compared to wild type by one way ANOVA). Scale bars correspond to 10  $\mu$ m (a), 100  $\mu$ m (c), and 50  $\mu$ m (d).



**FIGURE 5. Modulation of pathways downstream of ErbB signaling augments iPS-derived RBC generation.**

RNA-seq of BM CD34<sup>+</sup> cells treated with Lapatinib (increased affinities for EGFR and HER2), as well as Neratinib, Dacomitinib, and Afatinib (pan ErbB inhibitors) for 24 hours. Expression of genes from the Hallmark pathways (GSEA) significantly (\* = FDR<0.25) (a) upregulated and (b) downregulated (compared to DMSO vehicle) and common to pan-ErbB inhibitor treated cells. (c) Cell numbers (log fold change) of expanded iPS-derived CD34–5F cells<sup>35</sup> over 14 days. (d) Erythroid profile with and without CHIR99021 treatment during the final week of differentiation, with significantly increased proportion of mature (GlyA+, CD71-) erythrocytes. \* = p < 0.001 by unpaired, two-sided t-test. (e) morphology of May Grunwald Giemsa stained, iPS-derived orthochromatic erythroblasts. (c) n=7, (d) n=6 independent replicates from 2 experiments. Scale bar in (e) = 20 μm.