1   **Title:**

2   Gene-Centric Functional Dissection of Human Genetic Variation Uncovers

3   Regulators of Hematopoiesis

4

5   **Authors and Affiliations**

6   Satish K. Nandakumar[1,2,12], Sean K. McFarland[1,2,12], Laura Mateyka[1,2,3,13], Caleb A. Lareau[1,2,4,13],

7   Jacob C. Ulirsch[1,2,4,13], Leif S. Ludwig[1,2], Gaurav Agarwal[1,2,5,6], Jesse M. Engreitz[2,7], Bartlomiej

8   Przychodzen[8], Marie McConkey[9,] Glenn Cowley[2], John G. Doench[2], Jaroslaw P. Maciejewski[8],

9   Benjamin L. Ebert[2,9,10], David E. Root[2], Vijay G. Sankaran[1,2,6,11]

10

11   [1] Division of Hematology/Oncology, The Manton Center for Orphan Disease Research, Boston

12   Children's Hospital and Department of Pediatric Oncology, Dana-Farber Cancer Institute,

13   Harvard Medical School, Boston, MA, USA

14

15   [2] Broad Institute of MIT and Harvard, Cambridge, MA, USA

16

17   [3] Biochemistry Center (BZH), Ruprecht-Karls-University Heidelberg, 69120 Heidelberg,

18   Germany

19

20   [4] Program in Biological and Medical Sciences, Harvard Medical School, Boston, MA, USA

21

22   [5] University of Oxford, Oxford, UK

23

24   [6] Harvard Stem Cell Institute, Cambridge, MA, USA

25

26   [7] Harvard Society of Fellows, Harvard University, Cambridge, MA USA

27

28   [8] Department of Translational Hematology and Oncology Research, Taussig Cancer Institute,

29   Cleveland Clinic, Cleveland, OH, USA

30

31   [9] Division of Hematology, Brigham and Women's Hospital, Boston, MA, USA

32

33   [10] Department of Medical Oncology, Dana-Farber Cancer Institute, Boston, MA, USA

34

35   [11] Corresponding author

36

37   [12] Co-first authors

38

39   [13] Co-second authors

40

41

42   Correspondence: sankaran@broadinstitute.org

43   Keywords: hematopoiesis, genome-wide association studies, human genetics, functional

44   screen, erythropoiesis

45

46

## Abstract

48    Genome-wide association studies (GWAS) have identified thousands of variants associated with

49    human diseases and traits. However, the majority of GWAS-implicated variants are in non-

50    coding regions of the genome and require in depth follow-up to identify target genes and

51    decipher biological mechanisms. Here, rather than focusing on causal variants, we have

52    undertaken a pooled loss-of-function screen in primary hematopoietic cells to interrogate 389

53    candidate genes contained in 75 loci associated with red blood cell traits. Using this approach,

54    we identify 77 genes at 38 GWAS loci, with most loci harboring 1-2 candidate genes.

55    Importantly, the hit set was strongly enriched for genes validated through orthogonal genetic

56    approaches. Genes identified by this approach are enriched in specific and relevant biological

57    pathways, allowing regulators of human erythropoiesis and modifiers of blood diseases to be

58    defined. More generally, this functional screen provides a paradigm for gene-centric follow up of

59    GWAS for a variety of human diseases and traits.

60

61

## Introduction

As genotyping technologies and accompanying analytical capabilities have continued to improve, genome-wide association studies (GWAS) have identified tens of thousands of variants associated with numerous human diseases and traits. Despite these advances, our ability to discern the underlying biological mechanisms for the vast majority of such robust associations has remained limited, with a few exceptions (Claussnitzer et al., 2015; Gupta et al., 2017; Mohanan et al., 2018; Musunuru et al., 2010; Sankaran et al., 2008; Smemo et al., 2014). In general, published successes have required in-depth mechanistic studies of individual loci and implicated genes to decipher biological mechanisms.

Recent innovations in functional and computational genomics have advanced the field and enabled more rapid and higher-throughput identification of putative causal variants. Approaches that have shown the most success include the use of massively parallel reporter assays to examine allelic variation (Tewhey et al., 2016; Ulirsch et al., 2016; Vockley et al., 2015) and perturbation approaches for dissecting the necessity of regulatory elements (Fulco et al., 2016; Simeonov et al., 2017). In addition, genetic fine mapping approaches have improved our ability to identify putative causal variants among larger sets of variants in linkage disequilibrium (Guo et al., 2016; Huang et al., 2017; Lareau et al., 2018). However, even when putative causal variants are identified at a disease or trait-associated locus, they most often localize to non-coding regions of the genome, making it difficult to connect variants to genes that mediate the observed effects in a scalable manner (Claussnitzer et al., 2015; Gupta et al., 2017; Smemo et al., 2014).

85    In the context of hematopoiesis, GWAS studies have identified thousands of variants associated

86    with various blood cell traits, including hundreds associated with red blood cell traits alone (Astle

87    et al., 2016; van der Harst et al., 2012). Thorough follow-up efforts at individual loci have

88    identified important regulators of hematopoiesis, such as the key regulator of fetal hemoglobin

89    expression, BCL11A (Basak et al., 2015; Liu et al., 2018; Sankaran et al., 2008). However, as in

90    other tissues, the low-throughput with which associated genetic variants can be connected to

91    target genes underlying phenotypes continues to pose a problem for gaining biological insights

92    and clinical actionability in complex traits and diseases.

93

94    To accelerate the rate at which genetic variants can be connected to target genes, high-

95    throughput loss-of-function screens involving putative causal genes underlying the genetic

96    associations can be undertaken. This approach is complementary to conventional variant-

97    focused methods and overcomes bottlenecks that can arise during downstream target gene

98    identification. As a proof-of-principle, we connected variants associated with RBC traits to genes

99    regulating erythropoiesis by directly perturbing all candidate genes in primary human

100   hematopoietic stem and progenitor cells (HSPCs) undergoing synchronous differentiation into

101   the erythroid lineage. We demonstrate unique opportunities to rapidly implicate likely causal

102   genes and identify networks of biological actors underlying trait-associated variation. We

103   additionally illustrate the value of such screens to uncover previously unappreciated biological

104   regulators of human hematopoiesis that may serve as key disease modifiers.

105

106    # Results

107    **Design and Execution of an shRNA Screen Using Blood Cell Trait GWAS Hits to Identify**

108    **Genetic Actors in Erythropoiesis**

109    We applied a gene-centric loss-of-function screening approach to GWAS of RBC traits. We

110    focused on 75 loci associated with RBC traits that were identified by a GWAS performed in up to

111    135,000 individuals (van der Harst et al., 2012) spanning 6 RBC traits (Figure S1A). Importantly,

112    these 75 loci have been robustly replicated and show large effect sizes in more recently

113    reported association studies performed on larger cohorts and thus represent ideal targets for

114    perturbation studies (Astle et al., 2016; Lareau et al., 2018). We endeavored to select all genes

115    that could potentially underlie these 75 GWAS signals. To do this, each of the 75 sentinel SNPs

116    was first expanded to a linkage disequilibrium (LD) block including all SNPs in high LD ($r^2 > 0.8$,

117    Figure 1A, Figure S1B), then further to the nearest genomic recombination hotspot. Based upon

118    insights from previous expression quantitative trait locus (eQTL) studies (Montgomery and

119    Dermitzakis, 2011; Rossin et al., 2011; Veyrieras et al., 2008), each gene annotated in the

120    genome was expanded to include a wingspan encompassing 110 kb upstream and 40 kb

121    downstream of the transcriptional start and end sites, respectively, to also capture potential

122    functional regulatory elements. This resulted in selection of 389 genes overlapping or in the

123    vicinity of the LD blocks to be tested in the pooled loss-of-function screen. These were

124    distributed at a median of 4 genes per loci (Figure S1C).

125

126    Since the majority of common genetic variation underlying RBC traits appears to act in a cell-

127    intrinsic manner within the erythroid lineage, we decided to perturb the candidate genes during

128    the process of human erythropoiesis (Giani et al., 2016; Sankaran et al., 2012; Sankaran et al.,

129    2008; Ulirsch et al., 2016). We chose a pooled short hairpin RNA (shRNA) based loss-of-

130    function approach in primary hematopoietic cells to leverage a number of distinct strengths.

131    First, we have had prior success validating individual genes underlying RBC traits using shRNA-

132    based approaches in primary CD34$^+$ HSPC-derived erythroid cells. Second, shRNA libraries can

133    be much more efficiently packaged into lentiviruses and delivered to primary CD34$^+$ cells

134    compared to alternative CRISPR/Cas9-based guide RNA libraries (Ting et al., 2018). Third, it

135    avoids potential complications like non-uniform loss-of-function or gain-of-function outcomes

136    produced by CRISPR/Cas9 based approaches due to unpredictable DNA repair processes

137    (Mandegar et al., 2016). Furthermore, shRNAs can act rapidly to achieve gene knockdown and

138    thereby avoid compensatory effects that can occur when complete CRISPR knockout is

139    achieved (Rossi et al., 2015), better recapitulating the subtle changes in gene expression that

140    are characteristic of common genetic variation.

141

142    Mobilized peripheral blood-derived primary human CD34$^+$ HSPCs from 3 independent healthy

143    donors were infected with a lentiviral-based pooled shRNA library consisting of 2803 hairpins

144    targeting the 389 GWAS-nominated genes, along with 30 control genes (Moffat et al., 2006).

145    Each gene was targeted with 5-7 distinct shRNAs (Figure S1D). The set of control shRNAs

146    encompassed essential control genes, negative controls (e.g. luciferase and other non-

147    expressed genes), and a well-defined set of genes important for erythropoiesis (erythroid

148    controls) (Figure 1B, Table S2). Using lentiviral libraries with defined titers, we achieved an

149    infectivity of 35-50%, which is optimal for obtaining either zero or one stably integrated shRNA

150    per cell, while minimizing the possibility of infection of a single cell by multiple viruses. To

151    achieve sufficient library representation, we infected at least 1000 CD34$^+$ HSPCs per hairpin

152    (7~11 * 10$^6$ cells per experiment). The infected HSPCs were cultured using a three-phase

153    erythroid differentiation method (Giani et al., 2016; Hu et al., 2013) that results in synchronous

154    differentiation and maturation of erythroid progenitors into RBCs. We hypothesized that hairpins

155   targeting potential regulators of erythropoiesis would be depleted or enriched during the three-

156   phase erythroid culture, similar to our prior experience in analyzing specific GWAS-nominated

157   genes (Giani et al., 2016; Sankaran et al., 2012; Ulirsch et al., 2016). To assay these hairpins,

158   we isolated and deep-sequenced genomic DNA from the pool of infected cells at 6 different

159   culture time points that represent distinct stages of erythropoiesis to most broadly assess

160   putative causal genes that may act across the span of differentiation (Figure 1C, S1E).

161

162   **Summary Characterization of shRNA Screen Outcomes**

163   For the vast majority of the ~3000 hairpins included in the library, infection was efficient and

164   consistent. Greater than 95% of hairpins were represented at levels of at least 5 $\log_2$ counts per

165   million (CPM) at day 4, two days post-infection (Figure 2A). Across the two-week time course, a

166   diversity of effects - in terms of both increased and decreased hairpin abundance - were

167   observed. While many hairpins were selected against during the course of erythroid

168   differentiation, as reflected in decreases of those hairpin abundances over time, there were also

169   a number of hairpins that increased in the culture over the time course (Figure 2B, S2A).

170

171   The tested set of hairpins targeting genes nominated by the 75 loci showed a variety of

172   activities, forming a broad distribution spanning both decreases and increases in abundance at

173   different time points (Figure 2C). The various controls included in the library behaved as

174   expected. Hairpins targeting genes with known significance to erythropoiesis, such as *GATA1*

175   and *RPS19* (Khajuria et al., 2018; Ludwig et al., 2014), showed markedly decreased abundance

176   across the time course. Likewise, hairpins targeting a set of broadly essential genes (Table S2)

177   were strongly depleted by day 16 when compared to negative control hairpins targeting non-

178   human genes, which showed little if any change (Figure 2C-E). These trends were recapitulated

179   with strong correlation in each of the three donor CD34[+] cell backgrounds (Figure S2B).

180

**Statistical Modeling of Gene Effects and Accounting for Confounders in the shRNA**

**Screen**

The resulting longitudinal observations of hairpin abundance at each time point were used to

model the importance of each targeted gene during the process of erythropoiesis. A linear mixed

model was implemented to account for the longitudinal nature of the time course data (Li et al.,

2015) and to handle the confounding off-target and efficiency effects inherent to the shRNA

modality (Riba et al., 2017; Tsherniak et al., 2017). Since we wanted our model to be able to

detect significant changes in hairpin abundance at any time point throughout the differentiation

process, we converted the absolute hairpin abundances at each of the six time points to a $\log_2$

fold change relative to the initial hairpin abundances at the start of the differentiation. Using this

metric as our response variable, we specified a fixed effect for each gene to capture the

contribution that suppressing it with shRNAs would have on the respective abundances for each

of the resulting five time intervals. Given the potential variability that could emerge by using

shRNAs, we fit a random effect for each hairpin to minimize the chance of conflating inefficiency

or off-target effects with the specific on-target gene effect.

After fitting this model to the data, we selected our hit set using a two-threshold approach in

which both the magnitude and statistical confidence of the estimated gene effect size were

considered. Specifically, genes were called as hits if they had a fitted slope > 0.1 $\log_2$ fold

change per day within the interval while simultaneously possessing a Wald chi-square FDR-

adjusted q value < 0.1. This combined approach allowed us to avoid focusing on genes with

large, but highly variable or conflicted effects, as well as genes with highly confident but

miniscule effects. In total, this approach identified 77 genes at 38 of the 75 targeted loci which,

when suppressed, had a significant effect on the slope of shRNA-encoding DNA abundance at

205    any point during the time course. A majority of these hit loci (27 loci) had 1-2 gene targets

206    prioritized (Figure 3A, S3A). These candidate genes were found to be distributed across all 6 of

207    the originally annotated RBC GWAS traits (Figure S3B).

208

209    To evaluate the validity of this hit set, we began by assaying for enrichment of erythroid

210    essentiality, as recently quantified for each gene in the K562 erythroid cell line (Wang et al.,

211    2015). A permutation comparing the sum of K562 essentiality scores for the hit genes with those

212    of randomly drawn, identically-sized gene sets from the library of targeted genes revealed that

213    the hit set was indeed enriched with p = 0.0269 (Figure 3B). Likewise, when compared to

214    permuted sets of 77 genes randomly chosen from the genome (Figure S3C), there was even

215    stronger enrichment for erythroid essentiality with p = 0.00021, consistent with the idea that

216    genes in the library likely have stronger essentiality due to their genomic proximity to the GWAS

217    hits. We further explored whether the enrichment could be due to an intrinsic bias inherent to

218    GWAS screening itself by permuting sets of genes from libraries nominated by SNPs associated

219    with low-density lipoprotein levels, high-density lipoprotein levels, and triglyceride levels, finding

220    the hit set to be significantly enriched in all comparisons (Willer et al., 2013) (Figure S3D-F).

221

222    We further validated the ability of this approach to discover genetically relevant hits by

223    performing a permutation analysis based upon five "gold standard" genes in the library, which

224    possess known genetic underpinnings via identified causal variants: *CCND3* (Lareau et al.,

225    2018; Sankaran et al., 2012), *SH2B3 (Giani et al., 2016)*, *MYB* (Galarneau et al., 2010;

226    Sankaran et al., 2013; Sankaran et al., 2011), *KIT* (Jing et al., 2008; Lareau et al., 2018), and

227    *RBM38* (Ulirsch et al., 2016). Calculating the rank sums of hairpins ordered by our model's

228    computed FDR scores for 1,000,000 random combinations of five genes from the library yielded

229    a distribution over which enrichment for the five gold standards was seen with p=0.0249 (Figure

230   3C). While the vast majority of putative causal variants at the RBC trait-associated loci are in

231   non-coding regions, which can be challenging to use to identify a specific target gene, a subset

232   are in coding regions and thereby nominate a specific gene. As a result, we assayed for the

233   presence of coding variants fine-mapped to the interrogated loci from a recent large GWAS that

234   demonstrated a minimum posterior probability of association of 0.1 among the gene hits and

235   compared this with the overall set of genes interrogated in our library (Lareau et al., 2018).

236   Among the 389 GWAS-nominated genes in our library, 20 (~5%) were found to contain at least

237   one coding variant from this list. Of these, there was a significant enrichment observed among

238   the hits (~9%, p=0.03907 as determined by permutation analysis; Figure 3D).

239

240   Having established genetic confidence in our hit set, we next investigated whether the selected

241   genes satisfied enrichment criterion within the erythroid branch of hematopoiesis. RNA

242   expression values for each of the 77 hit genes were examined in datasets spanning human

243   hematopoiesis (Corces et al., 2016), as well as adult and fetal erythropoiesis (Yan et al., 2018)

244   (Figure 3E,F; Figure S3G). In the more holistic hematopoiesis dataset, common myeloid

245   progenitors (CMPs) and megakaryocyte-erythroid progenitors (MEPs) were significantly

246   enriched for hit genes (p < 0.01). These progenitor populations are known to contain the

247   progenitors that give rise to erythroid cells. Within a more detailed and separate analysis of

248   human adult erythropoiesis, proerythroblast, early basophilic, and late basophilic erythroblast

249   stages were particularly enriched (p < 0.001). The stage at which given genes are implicated to

250   play a role in erythropoiesis from the literature likewise often corresponded with the largest

251   magnitude fold changes across the longitudinal time course measurements, as was the case for

252   earlier genes like *RPL7A*, *RPL23A*, *RPS19*, and *KIT* (Gazda et al., 2012; Jing et al., 2008;

253   Moniz et al., 2012) as well as late genes like *SLC4A1* and *ANK1* (Bennett and Stenbuck, 1979;

254   Peters et al., 1996). Taken together, these results show that this functional gene-centric screen

255 allows for the identification of putative causal genes underlying RBC-trait GWAS hits, which

256 demonstrate clear enrichment in independent genetic and cell biological datasets. We are

257 therefore able to validate the utility of such an approach to identify biologically-relevant genes

258 underlying human genetic variation and holistically identify potential stages by which such target

259 genes may act to impact the process of hematopoiesis.

260

261 **Analysis of Interactions Among Members of the Hit Set Identifies Signaling, Structural,**

262 **and Translation-Related Subnetworks Important to Erythropoiesis**

263 By screening all loci and genes at once, our approach afforded us the immediate value of

264 examining mechanisms underlying the associations in a holistic fashion, unearthing both familiar

265 and more novel core gene network cassettes that play a role in erythropoiesis (Boyle et al.,

266 2017). Using STRING interaction network analyses (version 10.5) (Szklarczyk et al., 2017), we

267 could identify connectivity between the underlying nodes that highlighted a number of interacting

268 biological processes of both known and previously unappreciated importance to erythropoiesis

269 (Figure 4). We observed a number of molecules that play roles in cell signaling or transcriptional

270 regulation. MYB is a master regulator transcription factor that has been implicated in playing a

271 role in fetal hemoglobin regulation and in erythropoiesis more generally (Mucenski et al., 1991;

272 Wang et al., 2018). The *MYB* locus has been associated with numerous red blood cell traits

273 (including mean corpuscular volume, mean corpuscular hemoglobin concentration, and RBC

274 count) (Sankaran et al., 2013; van der Harst et al., 2012). ETO2 (CBFA2T3) is a part of the

275 erythroid transcription factor complex containing TAL1 and is required for expansion of erythroid

276 progenitors (Goardon et al., 2006). Both stem cell factor receptor KIT and erythropoietin

277 receptor (EPOR) mediated signaling are essential for erythropoiesis. Our screen identified KIT

278 as one of the factors underlying common genetic variation. CCND3 fills a critical role in

279 regulating the number of cell divisions during terminal erythropoiesis  and has been validated as

280    a causal gene associated with variation in RBC counts and size (Lareau et al., 2018; Sankaran

281    et al., 2012).

282

283    Interacting networks of hits also emerged in other aspects of red blood cell differentiation and

284    function. One of these centered around membrane and structural cytoskeletal proteins. Our

285    method recovered characteristic RBC genes like solute carrier family 4 member 1 (SLC4A1),

286    also known as band 3, (Peters et al., 1996), which serves as a key component of the RBC

287    membrane skeleton. Likewise, it recovered a direct interacting partner for SLC4A1, ankyrin 1

288    (ANK1), which anchors the cytoskeleton and cell membrane (Bennett and Stenbuck, 1979), as

289    well as N-ethylmaleimide Sensitive Factor, vesicle fusing ATPase (NSF), which facilitates

290    membrane vesicle trafficking within the cell (Glick and Rothman, 1987).

291

292    Within the realm of mRNA translation, a number of genes emerged as hits that specifically

293    highlight the role of the ribosome. This is interesting in light of recent work that has begun to

294    illuminate erythroid-specific effects of ribosomal perturbations (Khajuria et al., 2018; Ludwig et

295    al., 2014), although a connection between translation and common genetic variation affecting

296    RBC traits has not been previously appreciated. Both *RPL7A* and *RPL19*, for instance, have

297    been implicated by mutations observed in studies of Diamond-Blackfan anemia (Gazda et al.,

298    2012; Moniz et al., 2012). The common genetic variation affecting these ribosomal protein

299    genes might contribute to the incomplete penetrance and variable expressivity of anemia seen

300    in Diamond-Blackfan anemia patients (Ulirsch et al., 2018). Similar effects have been reported

301    in neurodevelopmental disorders, where common genetic variants may influence phenotypic

302    outcomes in patients (Niemi et al., 2018). Non-ribosomal hits in the mRNA metabolism space

303    were also found with both previously established and unknown ties to erythroid-specific

304    phenotypes. Exosome component 9 (EXOSC9), for instance, has been demonstrated previously

305    to act as part of the exosome complex as a specific gatekeeper of terminal erythroid maturation

306    (McIver et al., 2014). Other unappreciated components, including the tRNA methyltransferase

307    TRMT61A, also were highlighted through this analysis.

308

**309    Transferrin receptor 2 is a Negative Regulator of Human Erythropoiesis**

310    We selected several candidate genes identified by our screen for further validation, given their

311    previously unappreciated roles in human hematopoiesis/ erythropoiesis. The first, transferrin

312    receptor 2 (*TFR2*), encodes a protein canonically involved in iron homeostasis that has recently

313    been shown to also regulate EPO receptor signaling (Forejtnikova et al., 2010; Nai et al., 2015).

314    Although TFR2 has been studied in the context of murine erythropoiesis, its role in human

315    erythropoiesis has not been assessed. To validate TFR2 as a regulator of human erythropoiesis,

316    we performed individual knockdown experiments using lentiviral shRNAs in primary human

317    CD34$^+$ HSPCs undergoing erythroid differentiation. Significant knockdown of TFR2 was

318    observed at both the mRNA (Figure 5A) and protein levels (Figure 5B) using two independent

319    shRNAs from among the six targeting TFR2 in the screen. Though two of the six were outliers,

320    the two chosen here for follow-up were part of the consensus group of four showing similar

321    effects. Downregulation of TFR2 increased erythroid differentiation as observed by increased

322    expression of erythroid specific cell surface markers CD235a and CD71 at day 9 (shLUC ~22%;

323    TFR2 sh1 ~42%; TFR2 sh2 ~40%) and day 12 of culture (shLUC ~60%; TF2 sh1 ~80%; TFR2

324    sh2 ~80%) (Figure 5C, E & S5A). Downregulation of TFR2 also improved the later stages of

325    erythroid differentiation/ maturation, as observed by an increased rate of enucleation and

326    through assessment of cell morphology (Figure 5D & S5B). Previous studies have reported the

327    isolation of TFR2 as a component of the erythropoietin (EPO) receptor complex (Forejtnikova et

328    al., 2010). To test if downregulation of TFR2 can result in increased EPO signaling, we

329    measured EPO-dependent STAT5 phosphorylation after TFR2 knockdown in UT7/EPO cells

330    (Figure S5C). TFR2 downregulation resulted in significantly higher pSTAT5 phosphorylation in

331    comparison to the control with EPO stimulation from 0.02 U/mL to 200 U/mL (Figure 5F). In

332    addition, the maximal pSTAT5 response could be achieved within a shorter period of EPO

333    stimulation upon TFR2 downregulation (Figure S5D). Given our findings that TFR2 is a negative

334    regulator of EPO signaling, it may be an ideal therapeutic target for conditions characterized by

335    ineffective erythropoiesis like β-thalassemia (Rund and Rachmilewitz, 2005). A recent study has

336    supported this hypothesis, showing that Tfr2 downregulation is beneficial in a mouse model of

337    β-thalassemia (Artuso et al., 2018).

338

339    **SF3A2 is a Key Regulator of Human Erythropoiesis and is a Disease Modifier in a Murine**

340    **Model of Myelodysplastic Syndrome.**

341    Extensive mRNA splicing occurs during the terminal stages of erythropoiesis (Pimentel et al.,

342    2016). However, key regulators of this process remain largely undefined. Our study uncovered

343    splicing factor 3A subunit 2 (SF3A2) in the subnetwork of erythropoiesis signaling and

344    transcription hits (Figure 4).  SF3A2 specifically was associated with maximal hairpin drop out at

345    day 12 (FDR = 0.005) – a later time point in erythropoiesis. SF3A2 is a component of the

346    U2SNRP complex whose binding to the branch point is critical for proper mRNA splicing

347    (Gozani et al., 1996; Gozani et al., 1998). Knockdown of SF3A2 in primary human CD34[+]

348    HSPCs results in decreased cell numbers during erythroid differentiation starting from day 7

349    (Figure 6A-C). To measure early effects of SF3A2 and to exclude potential toxicity of puromycin

350    selection, we replaced the puromycin resistance gene with a GFP encoding cDNA in the

351    lentiviral shRNA constructs. We achieved similar infection (30~40% on day 6) at the early time

352    points between controls (shLuc) and shRNAs targeting *SF3A2* (Figure S6A). During erythroid

353    differentiation, we observed a reduction in GFP-expressing cells comparable to the decreased

354    cell numbers seen with the puromycin resistant constructs (Figure S6A). Decreased cell

355  numbers were associated with decreased erythroid differentiation as measured by erythroid

356  surface markers CD71 and CD235a (Figure 6D). We also observed an increase in non-erythroid

357  lineages based on surface marker expression of CD11b (myeloid) and CD41a (megakaryocyte)

358  (Figure S6B).

359

360  To identify the molecular mechanisms underlying the reduced differentiation of erythroid cells we

361  sorted stage-matched CD71$^+$/CD235$^+$ cells and performed RNA-Seq analysis. We also ran this

362  analysis in parallel for data from hematopoietic progenitors from patients with myelodysplastic

363  syndrome (MDS), a disorder well-known for significant impairment in terminal erythropoiesis,

364  either with or without somatic mutations in the related splicing factor *SF3B1 (Obeng et al.,*

365  *2016).* Cells treated with shRNA to suppress SF3A2 were found to differentially express 6061

366  genes with an adjusted p value < 0.05 as compared to the shLuc control, whereas only 807

367  genes were differentially expressed given the same threshold cutoff in the MDS patients with an

368  *SF3B1* mutations compared to those without (Figure 6E). Genes from both the SF3A2

369  differentially expressed set and the SF3B1 differentially expressed set were significantly

370  enriched for structural constituents of the ribosome (p < 3.2 x 10$^{-44}$ and p < 7.5 x 10$^{-24,}$

371  respectively) among other cellular components and functions (Tables S6-S9). Examining

372  differential splicing in the set of genes not differentially expressed in either condition, both were

373  found to exhibit a similar proportion of altered splicing events, including alternative 3' splice

374  sites, alternative 5' splice sites, mutally exclusive exons, and skipped exons with bayes factor >

375  10 (Figure 6E, Tables S10-S19).

376

377  We therefore wanted to further explore this connection between SF3A2 and its role in common

378  variation in RBC traits with SF3B1 and the role it plays in the pathogenesis of MDS. To this end,

379  we utilized a recently developed faithful mouse model harboring the *Sf3b1*$^{K700E}$ mutation that

380    displays characteristic features of MDS, including an anemia due to impaired erythropoiesis

381    (Obeng et al., 2016). We tested if downregulation of SF3A2 could worsen the already impaired

382    erythropoiesis seen in these animals. Equal numbers of lineage-negative HSPCs were isolated

383    from bone marrow of wild-type and *Sf3b1*[K700E] mice and infected with shRNAs targeting SF3A2

384    and then erythroid differentiation was induced (Figure S6C, D). Consistent with previous reports,

385    we observed that *Sf3b1*[K700E] cells show reduced erythroid differentiation and cell growth

386    compared to wild-type cells infected with control non-targeting shRNAs (Figure 6F, G, S6E-G).

387    Downregulation of SF3A2 using two independent shRNAs further worsens the defects in both

388    erythroid differentiation and cell growth observed for *Sf3b1*[K700E] cells (Figure 6F, G, S6E, F, G).

389    This data suggests that modulation of SF3A2 could modify the alterations of erythropoiesis

390    observed in the setting of somatic *SF3B1* MDS-causal mutations. This form of MDS is

391    characterized by significant variation in the degree of anemia found at the time of presentation

392    (Papaemmanuil et al., 2011). We therefore attempted to examine whether such common

393    genetic variation could contribute to such phenotypic variation. We identified a coding SNP,

394    rs25672, that was in LD with the sentinel SNP at this locus, rs2159213 ($r^2$ = 0.737675 in CEU

395    1000GENOMES phase 3). Prevalence of the alternate "G" allele (which is associated with the

396    prevalence of the "C" effect allele in the van der Harst et al. locus) has a suggestive correlation

397    with an increase in hemoglobin levels (Figure S6H) that was likely insignificant due to the limited

398    number of patients studied here. Unfortunately, larger cohorts in such a relatively rare disorder

399    could not be identified. However, these findings suggest that the subtle variation noted in

400    populations at the *SF3A2* locus may more profoundly cause variation among individuals with an

401    acquired blood disorder, such as MDS, illustrating the value of such a gene-centric study to

402    identify potential disease modifiers.

## Discussion

A major challenge in moving from GWAS-nominated variants to function is to identify potential target genes systematically. While many functional follow up approaches focus on causal variants, we reasoned that a gene-centered approach may be complementary to other emerging methods and represent a scalable approach for gaining broad insights into GWAS. To this end, we designed and executed a GWAS-informed high-throughput loss-of-function screen to identify key players in primary human HSPCs undergoing erythroid differentiation. Such dynamic *in vitro* systems afford a unique window through which to longitudinally screen, enabling unique insight to be gained into inherently non-stationary biological processes like erythropoiesis. The screen identified 77 gene hits at 38 of the original 75 loci used to design the library. Collectively, these hits had strongly amplified essentiality in erythroid cell lines, included a significant proportion of known, genetically-linked "gold standard" erythroid genes, and were enriched for red blood cell trait-associated coding variants orthogonally identified through genetic fine-mapping. From a holistic perspective, the network of interacting gene hits highlighted a number of high-level biological components and pathways important for erythropoiesis, including specific signaling and transcription factors, membrane and structural components, and components involved in mRNA translation.

Functional follow-up on *SF3A2* and *TFR2*, two gene hits identified in the screen, were fruitful in elucidating mechanistic ties between alteration in mRNA splicing and EPO signaling activity, respectively, to observed perturbation of erythroid phenotypes. In addition, our studies suggest that at least SF3A2, and potentially other regulators such as some implicated mRNA translation factors, may be key disease modifiers that alter the impaired erythropoiesis seen in diseases like MDS or Diamond-Blackfan anemia. These outcomes strongly recommend this screening

427    approach as a rapid means to access the genetic mediators underlying human erythroid

428    differentiation, with the potential to much more rapidly derive actionable biology from GWAS

429    studies. Moreover, since shRNA-based loss-of-function screens are readily accessible and offer

430    demonstrated compatibility with primary cell model systems, we believe this approach provides

431    a method that is portable and can be applied across a variety of lines of biological inquiry.

432

433    However, it is not a universal solution, and there are certainly a number of considerations that

434    must be kept in mind regarding the extent to which this type of assay can be adopted across

435    other diseases and traits. We acknowledge for it to be useful to a given research question, a

436    suitable system capable of modeling the trait/ disease of interest must first exist, and for many

437    cellular systems this is often challenging. Fortunately, this is a shortcoming that will diminish

438    over time as our understanding of human biology and our ability to faithfully recapitulate *in vivo*

439    microenvironments and processes improves, though this may be a distant prospect for

440    exquisitely complex tissues like the brain or for traits/ diseases that involve a larger number of

441    cell types/ interactions. Likewise, the use of shRNAs as the vehicle for perturbation carries with

442    it unique challenges, chief among them the proclivity of shRNA to exert confounding off-target

443    effects when compared to CRISPR-based methods. While this is true and unavoidable, the

444    inclusion of appropriate controls, both at the experimental level and in modeling off-target

445    contributors to observed phenotypic effects, provide an effective means to address this issue

446    (Tsherniak et al., 2017). We chose to perform our screen in primary hematopoietic cells and

447    thus were partially limited experimentally to the use of shRNA-based suppressive approaches.

448    Finally, evidence has recently been published that the targets of identified non-coding variants

449    are occasionally not within linkage disequilibrium blocks in the genome (Whalen and Pollard,

450    2018). This does not necessarily conflict with our results, since we identify hits at only 38 of 75

451    examined loci and provides an intriguing direction for further work that may elucidate how

452    genetic and epigenomic structural blocks in the human genome can provide complementary

453    information.

454

455    Our data shows that gene-centric screens are valuable for GWAS follow-up. They are not

456    limited to red cell traits and may be useful for other human traits/ diseases, as has begun to be

457    shown in disease-systems like Type II diabetes (Thomsen et al., 2016). Data from such screens

458    can be integrated with complementary insights gleaned from variant-centric screens. Ultimately

459    this could accelerate our understanding of human hematopoiesis and other biological

460    processes, and aid in the development of applicable therapies.

## Acknowledgments

461

462  We thank members of the Sankaran laboratory for valuable comments and suggestions on

463  these studies. This work was supported by the National Institutes of Health grants R01

464  DK103794 and R33 HL120791, as well as the New York Stem Cell Foundation (to V.G.S.).

465  V.G.S. is a New York Stem Cell Foundation-Robertson Investigator.

## Author Contributions

466

467  V.G.S., D.E.R., S.K.N., S.K.M., J.C.U., G.C. and J.G.D. designed the study. S.K.N performed

468  the pooled shRNA screen experiments in primary cells. G.C. and J.G.D prepared pooled shRNA

469  lentiviral libraries and performed downstream deep sequencing of barcodes. D.E.R. oversaw

470  shRNA library preparation and analysis of the screen. S.K.N., L.M., and G.A. performed all the

471  follow-up experiments. S.K.M., C.A.L, J.C.U., and J.M.E. performed all computational analyses.

472  L.S.L performed RNA sequencing on SF3A2 knockdown samples. M.M. and B.L.E. provided

473  samples from $Sf3b1^{K700E}$ knock-in mice. B.P. and J.P.M. analyzed and provided MDS patient

474  samples. S.K.N., S.K.M., and V.G.S. wrote the manuscript with input from all authors. V.G.S.

475  supervised all experimental and analytic aspects of this work.

476

## Competing Interests:

477

478  The authors declare no competing interests.

## Figure Legends

479

480   **Figure 1: Design and Execution of an shRNA Screen Using Blood Cell Trait GWAS Hits to**

481   **Identify Genetic Actors in Erythropoiesis** (A) Overview of shRNA library design. 75 loci

482   associated with red blood cell traits (van der Harst et al., 2012) were used as the basis to

483   calculate 75 genomic windows of LD 0.8 or greater from the sentinel SNP. Genes with a start

484   site within 110 kb or end site within 40 kb of the LD-defined genomic windows were chosen as

485   candidates to target in the screen. (B) Compositional makeup of the library, depicted as number

486   of genes and number of hairpins for each of the four included subcategories; GWAS-nominated

487   genes, erythroid genes, essential genes, and negative control genes. (C) Primary CD34$^+$

488   hematopoietic stem and progenitor cells (HSPCs) isolated from 3 independent donors were

489   cultured for a period of 16 days in erythroid differentiation conditions. At day 2, cells were

490   infected with the shRNA library, and the abundances of each shRNA were measured at days 4,

491   6, 9, 12, 14, and 16 using deep sequencing.

492

493   **Figure 2: Summary Characterization of shRNA Screen Outcomes** (A) Kernel density plot

494   showing library representation as $\log_2$ shRNA CPM across all hairpins. (B) shRNA abundance

495   $\log_2$ fold changes from day 4 to day 16. Represented values are the mean of hairpin abundance

496   $\log_2$ fold changes across hairpins for each gene and two standard deviations. (C) Kernel density

497   plots representing the day 4 to day 16 $\log_2$ fold changes of hairpin abundances for each of the

498   subcategories of the library, including GWAS-nominated genes, known erythroid essential

499   genes, essential genes to cell viability, and orthogonal genes serving as negative controls. (D)

500   Violin plot of day 4 and day 16 $\log_2$ CPM for known actors *GATA1* and *RPS19* and negative

501   controls LacZ and luciferase. (E) $\log_2$ hairpin counts averaged for known actors *GATA1* and

502   *RPS19* as well as negative controls LacZ and luciferase across the course of the experiment.

503   Gray lines depict the universe of all other gene traces in the library for context.

504

505 **Figure 3: Statistical Modeling of Gene Effect Accounting for Off-target shRNA**

506 **Confounders** (A) Bar graph showing the 38 of 75 loci in the screen with at least one

507 corresponding statistically significant (FDR < 0.1, β > 0.1) gene effect causing either a positive

508 or negative log$_2$ fold change in shRNA abundance. (B) Kernel density plot showing the expected

509 distributions of K562 essentiality scores using permuted gene hit sets from the library. (C)

510 Hairpin rank sums for permuted sets of 5 genes. The red line indicates the enriched rank sums

511 for 5 "gold standard" genes included in the library, *CCND3*, *SH2B3*, *MYB*, *KIT*, and *RBM38*, for

512 each which a genetic basis of action has already been established. (D) Permuted distribution of

513 % inclusion of predicted coding variants among the set of identified hits. (E) Heat map depicting

514 strength of expression (as z scores within each gene) for each of the 77 identified hit genes

515 across hematopoietic lineages (top) and throughout the specific stages of adult erythropoiesis

516 (bottom). Purple boxes highlight the cell types that were enriched for expression of hit genes.

517 (F) Calculated enrichment of the identified hit genes for expression across hematopoietic

518 lineages (top) and throughout the specific stages of adult erythropoiesis (bottom). In both cases,

519 cellular states corresponding to those along the erythropoietic lineage had elevated probability

520 of expressing genes from the hit set as compared to other genes from the library.

521

522 **Figure 4: Analysis of interactions among members of the hit set identifies**

523 **signaling/transcription, membrane, and mRNA translation-related subnetworks important**

524 **to erythropoiesis.** Comparison of gene hits via the STRING database identified interacting

525 networks related to hematopoietic signaling/transcription, membrane, and mRNA translation-

526 related subnetworks**.** Edges of the network are color-coded according to the evidence

527 supporting the interaction.

528

529    **Figure 5: Transferrin receptor 2 is a Negative Regulator of Human Erythropoiesis** (A)

530    Quantitative RT-PCR and (B) Western blot showing the expression of TFR2 in human CD34$^+$

531    cells five days post-infection with the respective lentiviral shRNAs targeting TFR2 (TFR2 sh1 &

532    sh2) and a control luciferase gene (shLUC). (C) Representative FACS plots of erythroid cell

533    surface markers CD71 (transferrin receptor) and CD235a (Glycophorin A) expression at various

534    time points during erythroid differentiation. Percentages in each quadrant is represented as

535    mean and standard deviation of 3 independent experiments (D) Hoechst staining showing more

536    enucleated cells after TFR2 knockdown at day 21 of erythroid culture. (E) Representative

537    histogram plots showing increased expression of CD235a (Glycophorin A) after TFR2

538    knockdown (F) Enhanced pSTAT5 response after TFR2 knockdown in UT7/EPO cells.

539

540    **Figure 6: SF3A2 is a Key regulator of Human Erythropoiesis and Modulates**

541    **Erythropoiesis Defects in a Murine Model of MDS** (A) Quantitative RT-PCR and (B) Western

542    blot showing the expression of SF3A2 in human CD34$^+$ cells five days post-infection with the

543    respective lentiviral shRNAs targeting *SF3A2* (sh1 & sh2) and a control luciferase gene

544    (shLUC). (C)  Growth curves showing that downregulation of SF3A2 results in reduced total cell

545    numbers during erythroid differentiation from 3 independent experiments. (D) Representative

546    FACS plots of erythroid cell surface markers CD71 (transferrin receptor) and CD235a

547    (Glycophorin A) expression at various time points during erythroid differentiation. Percentages in

548    each quadrant is represented as mean and standard deviation of 3 independent experiments

549    (E) Altered splicing events identified by RNA-Seq analysis of stage matched erythroid cells

550    (shSF3A2 vs. shLUC). Overlapping changes observed in *SF3B1* mutant BM cells from MDS

551    patients (Obeng et al). (F) Lineage negative bone marrow cells from wildtype (WT) and

552    *Sf3b1$^{K700E}$* mice were infected with shRNAs targeting murine *Sf3a2* gene co-expressing a

553    reporter GFP gene. Percentage of Ter119$^+$ CD71$^+$ erythroid cells within the GFP compartment

554    after 48hrs in erythroid differentiation. (G) Total cell numbers of GFP$^+$ erythroid cells after 48hrs

555    in erythroid differentiation.

556

## Materials and Methods:

557

558  Design of the shRNA Library

559  Ensembl assembly GRCh37p9 was utilized to expand 75 SNPs previously identified in a RBC

560  trait GWAS to include a genomic region in linkage disequilibrium with $r^2 \geq 0.8$. Each of these

561  regions was then further expanded to the nearest recombination hotspot. All genes in the

562  genome were expanded to include 110 kb upstream and 40 kb downstream of the transcription

563  start and end sites, respectively, to maximize capture of non-coding regulatory interactions

564  based upon previously published observations. Genes with windows calculated in this way

565  found to be overlapping with any of the SNP windows were flagged for inclusion in the screen.

566  In addition, each locus was examined individually, and in cases of gene deserts, unusually

567  proximal recombination hotspots, or other unusual genomic structures, the SNP region was

568  expanded to include additional genes nearby. This resulted in a total of 389 test genes, which

569  were each targeted by 4-7 distinct shRNAs. Also included in the library were shRNAs targeting a

570  set of 8 validated erythroid genes (*GATA1, RPL5, RPS19, EPOR, ALAS2, CDAN1, SEC23B,*

571  *ZFPM1*). A pooled library of 2803 TRC clones was produced from the sequence-validated TRC

572  shRNA library (Moffat et al., 2006) and included shRNAs targeting control genes and essential

573  genes.

574

575  Pooled shRNA Screening

576  Mobilized peripheral blood CD34$^+$ cells from three separate donors (7~11 * 10$^6$ cells per donor)

577  were differentiated into erythroid cells using a three-stage system that has been previously

578  described (Hu et al., 2013). Cells were cultured using IMDM containing 2% human plasma, 3%

579    human AB serum, 200 μg/ml human holo-transferrin, 3 IU/mL heparin, and 10 mg/mL insulin

580    (base medium). During days 0 to 7, cells were supplemented with IL-3 (1 ng/mL), SCF (10

581    ng/ml), and EPO (3 IU/ml). On day 2 of this culture, cells were transduced with the pooled

582    lentiviral shRNA library prepared by Broad Institute Genetic Perturbation Platform (1ml of virus

583    per 0.75 * $10^6$ cells) by spinfection at 2000 rpm for 90 minutes with 6 μg/ml polybrene. During

584    days 7 to 13, cells were supplemented with SCF and EPO only. After day 13, cells were

585    supplemented with EPO alone and the holo-transferrin concentration was increased to 1 mg/ml.

586    A minimum of 10 * $10^6$ cells was re-plated at each time point to ensure appropriate library

587    representation and prevent bottlenecks among the infected cells. Cell pellets were made from

588    20~80 * $10^6$ cells at days 4, 6, 9, 12, 14, and 16. At the conclusion of the pooled screen,

589    genomic DNA (gDNA) was extracted from the cell pellets using NucleoSpin Blood XL-Maxi kit

590    (Clonetech) according to kit specifications. The shRNA-containing region was PCR amplified

591    from the purified gDNA and barcoded using the following conditions: 0.5 μl P5 primer mix

592    (100μM), 10 μl P7 primer mix (5μM), 8 μl dNTP mix, 1x ExTaq buffer, 1.5 μl of ExTaq DNA

593    polymerase (Takara), and up to 10 μg genomic DNA in a total reaction volume of 100 μl. A total

594    of 40~87.5 μg gDNA was used as template from each conditions. Thermal cycler PCR

595    conditions consisted of heating samples to 95 °C for 5 min; 28 cycles of 95 °C for 30 s, 53 °C for

596    30 s, and 72 °C for 20 s; and 72 °C for 10 min. Equal amounts of samples were then mixed and

597    purified using AMPure XP for PCR purification (Beckman Coulter). Samples were sequenced

598    using a custom sequencing primer using standard Illumina conditions by the Broad Institute

599    Genetic Perturbation Platform. Sequencing reads were deconvolved and hairpin counts were

600    quantified for subsequent analysis by counting against the barcode reference using PoolQ

601    (https://portals.broadinstitute.org/gpp/public/dir/download?dirpath=software&filename=poolq-

602    2.2.0-manual.pdf).

603

26

604     P5 primer

605     **AATGATACGGCGACCACCGAGATCT**ACACTCTTTCCCTACACGACGCTCTTCCGATCT[s]**TC**

606     **TTGTGGAAAGG*A*C*G*A**

607     A mix of P5 primers with stagger regions [s] of different length was used to maintain

608     sequence diversity across the flow-cell.

609

610     P7 primer

611     **CAAGCAGAAGACGGCATACGAGAT**NNNNNNNNGTGACTGGAGTTCAGACGTGTGCTCTTC

612     CGATCT**TCTACTATTCTTTCCCCTGCA*C*T*G*T**

613     Independently barcoded P7 primers was used for each condition.

614     NNNNNNNN – barcode region

615

616     <u>Analysis of the shRNA Screen</u>

617     Three separate donor primary CD34[+] cells populations were run as replicates in the shRNA

618     screen. A pseudocount of 1 was added to all shRNA-encoding DNA count totals and these

619     counts were subsequently normalized to counts per million (CPM) and $\log_2$ transformed. A linear

620     mixed model was constructed to fit fixed effects for each gene ($g$) using the $\log_2$ fold change

621     from initial hairpin counts as the response variable ($y$). A random effect was included to capture

622     variations in efficacy and off-target effects for each shRNA ($h$) used to target a given gene. The

623     resulting model, $y \sim g + (0 + h|g)$, was fit in R-3.4 using the lme4 package. Genes hits were

624     called from the set of genes with β coefficient effect size > 0.1 and the Wald chi-square test

625     adjusted q value < 0.1. Enrichment of erythroid essential genes within the hit set was calculated

626     by running 1 million permutations against the distribution of K562 essentiality for all genes

627     included in the library, panels of genes nominated by sets of significant GWAS-associated lipid

628     trait SNPs (Willer et al., 2013), and against all genes in the genome (Ensembl GRCh37p9).

27

629   Enrichment for identification of the included 5 "gold standard" genes and for red blood cell trait-

630   associated coding variants were each accomplished using identical permutation schemes.

631   Expression of the hit genes in various cell states/ stages of differentiation were derived from the

632   cited datasets and permuted across all unique stages to determine stage-specific enrichment.

633   The interaction network surrounding the 77 hits identified in the screen was generated in the

634   latest version of STRING (10.5) and filtered for the purposes of display to only those nodes with

635   at least one edge to another node among the hits.

636

637   RNA-Seq

638   Stage matched CD71$^+$/ CD235a$^+$ cells derived from CD34$^+$ HSPCs infected with SF3A2 sh3, sh4

639   and shLUC were FACS sorted at day 8 of erythroid differentiation. RNA was isolated using a

640   RNAqueous Micro kit (Invitrogen) according to the manufacturer's instructions. DNase digestion

641   was performed before RNA was quantified using a Qubit RNA HS Assay kit (Invitrogen). 1-10 ng

642   of RNA were used as input to a modified SMART-seq2 (Picelli et al., 2014) protocol and after

643   reverse transcription, 8-9 cycles of PCR were used to amplify transcriptome library. Quality of

644   whole transcriptome libraries was validated using a High Sensitivity DNA Chip run on a

645   Bioanalyzer 2100 system (Agilent), followed by library preparation using the Nextera XT kit

646   (Illumina) and custom index primers according to the manufacturer's instructions. Final libraries

647   were quantified using a Qubit dsDNA HS Assay kit (Invitrogen) and a high sensitivity DNA chip

648   run on a Bioanalyzer 2100 system (Agilent). All libraries were sequenced using Nextseq High

649   Output Cartridge kits and a Nextseq 500 sequencer (Illumina). Libraries were sequenced using

650   2x38bp paired end reads.

651

652   RNA-seq Differential Expression Analysis

653  For differential expression analysis, paired end sequencing reads from our SF3A2 shRNA

654  knockdown experiments and obtained from the SF3B1 mutant datasets (Obeng et al., 2016)

655  were quantified using Salmon version 0.11.1 (Patro et al., 2017) with default parameters and an

656  index constructed from Gencode annotations version 28. Differential expression of quantified

657  counts were calculated using DESeq2 (Love et al., 2014) in R-3.4. Enrichment for functions and

658  components of the cell among the differentially expressed gene sets were quantified using

659  GOrilla (Eden et al., 2007; Eden et al., 2009).

660

661  RNA-seq Differential Splicing Analysis

662  Paired end sequencing reads from our SF3A2 shRNA knockdown experiments and obtained

663  from the cited SF3B1 mutant datasets were aligned using STAR version 2.5.2 in two-pass

664  mode. Differential splicing was quantified using MISO version 0.5.4 in Python 2.7 using the

665  instructions and annotation files provided with the package (Katz et al., 2010).

666

667  Analysis of Hemoglobin Levels for MDS Patients with or without SF3A2 Mutations

668  Genotyped MDS patient hemoglobin level measurements were obtained from the laboratory of

669  J. Maciejewski. 1000GENOMES phase 3 data was used to find a SNP encoded in whole-exome

670  sequencing data (rs25672) in high LD ($r^2$ = 0.737675) with the SF3A2-associated sentinel SNP

671  (rs2159213). An ordinary least squares linear regression was used to fit the patient hemoglobin

672  levels to the number of SF3A2 minor alleles present in each patient (log likelihood ratio test p =

673  0.140).

674

675  Phosphorylated STAT5 assessment with Intracellular Flow Cytometry

676  UT-7/EPO cells were cultured in DMEM medium supplemented with 10% Fetal Bovine Serum

677  and 2 U/mL EPO. 5 days post-infection with TFR2 shRNAs, UT-7/EPO cells were cytokine

678    starved overnight. On the next day, cells were treated with EPO in a dose dependent manner

679    ((0 U/mL, 0.002 U/mL, 0.02 U/mL, 0.2 U/mL, 2 U/mL, 20 U/mL and 200 U/mL) and incubated 37

680    °C for 30 min. Alternatively the cells were treated with 2U/ml EPO in a time dependent manner

681    (15, 30, 60, 120,180 min). Treated cells were gently mixed with pre-warmed Fixation Buffer (BD

682    Bioscience) at 37°C for 10 min to fix cells. To permeabilize cells for intracellular staining, cells

683    were resuspended in pre-chilled Perm Buffer III (BD Bioscience) for 30 min at 4°C. After three

684    washes with 3% FBS in PBS, samples were stained either with Alexa Fluor-647 Mouse Anti-

685    phospho-STAT5 (pY694; 1:20 dilution) for 1 hr in the dark at room temperature. A BD Accuri C6

686    Cytometer (BD Bioscience) was used to acquire mean fluorescent intensity (MFI) of phospho-

687    STAT5-Alexa Fluor 647. The MFI of phospho-STAT5-Alexa Fluor 647 of gated single cells was

688    calculated using FlowJo (version 10.0.8r1). Unstimulated UT7/EPO cells were used as a

689    negative control.

690

691    <u>May-Grünwald-Giemsa staining</u>

692    Approximately 50,000 – 200,000 cells were harvested, washed once at 300 x g for 5 min,

693    resuspended in 200 µL FACS buffer and spun onto poly-L-lysine coated glass slides (Sigma

694    Aldrich) with a Shandon 4 (Thermo Fisher) cytocentrifuge at 300 rpm for 4 min. Visibly dry

695    slides were stained with May-Grünwald solution for 5 min, rinsed 4 times for each 30 s in $H_2O$,

696    transferred to Giemsa solution for 15 min and washed as described above. Slides were dried

697    overnight and mounted with coverslip. All images were taken with AxioVision software (Zeiss)

698    at 100 x magnification.

699

700    <u>Mouse erythroid differentiation culture</u>

701    Bone marrow cells were isolated from $SF3B1^{K700E\ +/-}$ mice and littermate controls were lineage

702    depleted using Lineage Cell Depletion Kit, mouse (Miltenyi Biotech) according to manufacturer's

703    protocols. Lineage negative cells were immediately transduced with lentiviral shRNAs targeting

704    SF3A2 or controls (MOI -90) by spinfection at 2000rpm for 90 min. The cells were cultured in

705    erythroid maintenance medium (StemSpan-SFEM; StemCell Technologies) supplemented with

706    100 ng/mL recombinant mouse stem cell factor (SCF) (R&D Systems), 40 ng/mL recombinant

707    mouse IGF1 (R&D Systems), 100 nM dexamethasone (Sigma), and 2 U/mL erythropoietin

708    (Amgen) and cultured at 37°C for 36 hours. Following this the cells were cultured for another 48

709    hours in erythroid differentiation medium (Iscove modified Dulbecco's medium containing 15%

710    (vol/vol) FBS (Stemcell), 1% detoxified BSA (Stemcell), 500 µg/mL holo-transferrin (Sigma-

711    Aldrich), 0.5 U/mL Epoetin (Epo; Amgen), 10 µg/mL recombinant human insulin (Sigma-Aldrich),

712    and 2 mM L-glutamine (Invitrogen)) at 37 °C.

713

714    <u>Flow cytometry analyses and antibodies</u>

715    All flow cytometry data was acquired using either using LSR II SORP or LSR Fortessa flow

716    cytometers (BD Biosciences). All staining was carried out in FACS buffer (2% FBS in PBS) for

717    30 minutes on ice unless otherwise described. The following antibodies were used anti-human

718    CD235a-APC (eBioscience, Clone HIR2), anti-human CD71-FITC (eBioscience, Clone OKT9),

719    anti-human CD71-PEcy7 (eBioscience, Clone OKT9), ant-human CD49d-PE (Miltenyi, Clone

720    MZ18-24A9), anti-human CD41a-PE (eBioscience, Clone HIP8), anti-human CD11b-PE

721    (eBioscience, Clone ICRF44), anti-mouse Ter119-APC (eBioscience, Clone TER119), anti-

722    mouse CD71-PE (eBioscience, Clone R17217) and Alexa Fluor-647 anti-phospho STAT5

723    (pY694) (BD Bioscience Cat#: 612599). Hoechst 33342 (Life Technologies, H1399) was used to

724    visualize nuclei.

725

726    <u>shRNA sequences</u>

727    The following lentiviral shRNA constructs were generated in Polymerase III based shRNA

728    backbone pLKO.1-puro (Sigma Aldrich).

729    shLUC

730    5'-CCGGCGCTGAGTACTTCGAAATGTCCTCGAGGACATTTCGAAGTACTCAGCGTTTTTG-3'

731    TFR2 sh1

732    5'-CCGGGCCAGATCACTACGTTGTCATCTCGAGATGACAACGTAGTGATCTGGCTTTTTG-3

733    TFR2 sh2

734    5'-CCGGCAACAACATCTTCGGCTGCATCTCGAGATGCAGCCGAAGATGTTGTTGTTTTTG-3'

735    SF3A2 sh1 (human)

736    5'-CCGGCTACGAGACCATTGCCTTCAACTCGAGTTGAAGGCAATGGTCTCGTAGTTTTT-3

737    SF3A2 sh2 (human)

738    5'-CCGGCCTGGGCTCCTATGAATGCAACTCGAGTTGCATTCATAGGAGCCCAGGTTTTT-3'

739    SF3A2 sh3 (human)

740    5'-CCGGCAAAGTGACCAAGCAGAGAGACTCGAGTCTCTCTGCTTGGTCACTTTGTTTTT-3

741    SF3A2 sh4 (human)

742    5'-CCGGACATCAACAAGGACCCGTACTCTCGAGAGTACGGGTCCTTGTTGATGTTTTTT-3'

743

744    The following lentiviral shRNA constructs were generated in Polymerase II based mir30 shRNA

745    backbone developed in the lab SFFV-Venus-mir30 shRNA backbone.

746    shNT(non-targeting)

747    5'_TGCTGTTGACAGTGAGCGATCTCGCTTGGGCGAGAGTAAGTAGTGAAGCCACAGATGTA
748    CTTACTCTCGCCCAAGCGAGAGTGCCTACTGCCTCGGA_3'
749

750    *Sf3a2* sh1 (mouse)

751    5'_TGCTGTTGACAGTGAGCGCGGAGGTGAAGAAGTTTGTGAATAGTGAAGCCACAGATGTA
752    TTCACAAACTTCTTCACCTCCATGCCTACTGCCTCGGA_3'
753
754    *Sf3a2* sh2 (mouse)

755  5'_TGCTGTTGACAGTGAGCGACCACCGTTTCATGTCTGCTTATAGTGAAGCCACAGATGTAT
756  AAGCAGACATGAAACGGTGGCTGCCTACTGCCTCGGA_3'

757
758  *Sf3a2* sh3 (mouse)

759  5'_TGCTGTTGACAGTGAGCGATCCTGCCTTGAGCCTATTAAATAGTGAAGCCACAGATGTAT
760  TTAATAGGCTCAAGGCAGGACTGCCTACTGCCTCGGA_3'

761
762  *Sf3a2* sh4 (mouse)

763  5'_TGCTGTTGACAGTGAGCGACCACTGGAACAGAGAAACCAATAGTGAAGCCACAGATGTA
764  TTGGTTTCTCTGTTCCAGTGGGTGCCTACTGCCTCGGA_3'

765
766  *Sf3a2* sh5 (mouse)

767  5'_TGCTGTTGACAGTGAGCGATGGAGGTGAAGAAGTTTGTGATAGTGAAGCCACAGATGTA
768  TCACAAACTTCTTCACCTCCACTGCCTACTGCCTCGGA_3'

769
770
771  qPCR primers

772  *TFR2* Fwd: 5'-ATCCTTCCCTCTTCCCTCCC-3'

773  *TFR2* Rev: 5'-CCATCCAGCCACATGGTTCT-3

774  *SF3A2* Fwd: 5'-CCTGAGAAGGTCAAGGTGGA-3'

775  *SF3A2* Rev: 5'-CTCCGAGTCTCTCTGCTTGG-3'

776

777  Western Blot antibodies

778  Anti-GAPDH (Santa Cruz Biotechnology, sc-32233); anti-TFR2 (Santa Cruz Biotechnology, sc-

779  sc-32271); anti-SF3A2 (Santa Cruz Biotechnology, sc-390444)

780

781

782

## 783 Supplemental Information

## 784 Supplementary Figure Legends

785 **Supplementary Figure S1: Summary Characteristics of Designed shRNA Library** (A)

786 Counts of loci from among the original 75 annotated with linkage to each of the six RBC traits,

787 hemoglobin (Hb), mean corpuscular hemoglobin (MCH), mean corpuscular hemoglobin

788 concentration (MCHC), mean corpuscular volume (MCV), packed cell volume (PCV), and red

789 blood cell count (RBC). Some loci were associated with multiple traits. (B) Kernel density plot

790 showing the $\log_{10}$ sizes in bp of the LD-defined genomic windows used to find overlapping

791 genes. (C) Histogram showing distribution of number of genes selected using the LD window

792 method at each locus. A median of 4 genes were present at each. (D) Histogram showing

793 distribution of number of independent hairpins included in the library to target each of the

794 candidate's genes. (E)  Representative FACS plots of erythroid cell surface markers CD71

795 (transferrin receptor) and CD235a (Glycophorin A) expression at various time points during

796 erythroid differentiation in uninfected (Mock) or CD34$^+$ cells infected with the shRNA library

797 (Pool). Percentages in each quadrant is represented as mean and standard deviation of 3

798 experiments from independent donors.

799

800 **Supplementary Figure S2: Additional Metrics of Library Performance** (A) shRNA

801 abundance $\log_2$ fold changes from day 4 to each of the other time points. Represented values

802 are the mean of hairpin abundance $\log_2$ fold changes across hairpins for each gene and two

803 standard deviations. (B) Scatter plots showing agreement of replicate observations across

804 independent CD34$^+$ donor populations.

805

34

806    **Supplementary Figure S3: Additional Characterization of Modeling Outcomes** (A)

807    Histogram showing the number of gene hits identified at each of the 40 loci with at least one

808    significant gene effect detected. (B) Bar graph showing the number of gene hits identified for

809    each of the 6 red blood cell traits used in the original GWAS to identify the studied loci. (C)

810    Permuted enrichment of essentiality among the set of hit genes vs. randomly chosen sets of

811    genes from the human genome. (D) Permuted enrichment of essentiality among the set of hit

812    genes vs. genes implicated by a separate GWAS for LDL cholesterol levels. (E) Permuted

813    enrichment of essentiality among the set of hit genes vs. genes implicated by a separate GWAS

814    for HDL cholesterol levels. (F) Permuted enrichment of essentiality among the set of hit genes

815    vs. genes implicated by a separate GWAS for blood triglyceride levels. (G) Heat map depicting

816    strength of expression (as z scores within each gene) for each of the 77 identified hit genes

817    throughout the specific stages of fetal erythropoiesis. Purple boxes highlight the cell types that

818    were enriched for expression of hit genes.

819

820    **Supplementary Figure S5: Transferrin Receptor 2 is a Negative Regulator of Human**

821    **Erythropoiesis** (A) Representative FACS plots of alternate erythroid cell surface markers

822    CD49d (α4 integrin) and CD235a (Glycophorin A) expression at various time points during

823    erythroid differentiation. (B) May-Grunwald Giemsa staining showing more differentiated

824    erythroid cells after TFR2 knockdown at day 18 of erythroid culture. (C) Western blot showing

825    downregulation of TFR2 in UT7/EPO cells. (D) Time dependent absolute value of MFI of STAT5

826    in UT7/ Epo cells after TFR2 knockdown.

827

828    **Supplementary Figure S6: SF3A2 is Required for Human Erythropoiesis and Modulates**

829    **Erythropoiesis Defects in a Murine Model of MDS** (A) shRNAs targeting *SF3A2* co-

830    expressing a reporter GFP gene was infected into CD34[+] cells and cultured in erythroid

831   conditions. GFP expression at various time points from three independent experiments show

832   that downregulation of SF3A2 results in reduced cell numbers. (B) Representative FACS plots

833   of erythroid (CD235a) and non-erythroid cell surface markers (CD11b / CD41a) and at various

834   time points showing an increase in non-erythroid lineages upon SF3A2 downregulation. Cells

835   were gated on the GFP positive population. (C) Knockdown efficiency of shRNAs targeting

836   *SF3A2* in murine erythroleukemia (MEL) cells by western blot. (D) Total cell numbers of GFP+

837   cells at the start of murine erythroid differentiation. (E) Percentage of Ter119$^+$ CD71$^+$ erythroid

838   cells within GFP compartment and (F) Total cell numbers of GFP$^+$ erythroid cells after 24hrs in

839   erythroid differentiation. (G) Growth curves of GFP$^+$ erythroid cells during erythroid culture. (H)

840   Putative but insignificant interaction between *SF3A2* variant alleles (rs25672) and hemoglobin

841   levels in MDS patients with *SF3B1* mutations.

842

843 ## Supplemental Tables

844 **Table S1:** Table containing annotations and information for the 75 SNPs used to seed the

845 shRNA library.

846 **Table S2:** Table containing annotations and information for all hairpins, as well

847 as shRNA counts for each time point and replicate.

848 **Table S3:** Table containing the R model output for each gene.

849 **Table S4:** Table containing the DESeq2 output for differentially expressed genes in cells

850 undergoing SF3A2 knockdown or control shRNA treatment.

851 **Table S5:** Table containing the DESeq2 output for differentially expressed genes in MDS

852 patients with and without mutations in *SF3B1*.

853 **Tables S6, S7:** Tables containing the GO component (S6) and function (S7) enrichments

854 calculated using GOrilla for cells undergoing SF3A2 knockdown or control shRNA treatment.

855 **Tables S8, S9:** Tables containing the GO component (S8) and function (S9) enrichments

856 calculated using GOrilla for MDS patient samples with and without mutations in *SF3B1*.

857 **Tables S10-S14:** Tables containing the differential splicing analysis for cells undergoing SF3A2

858 knockdown or control shRNA treatment. Categories of splice mutations presented in each table

859 are alternative 3' splice sites, alternative 5' splice sites, mutually exclusive exons, retrained

860 introns, and skipped exons, respectively.

861 **Tables S15-S19:** Tables containing the differential splicing analysis for MDS patient patient

862 samples with and without mutations in *SF3B1*. Categories of splice mutations presented in each

863 table are alternative 3' splice sites, alternative 5' splice sites, mutually exclusive exons, retrained

864 introns, and skipped exons, respectively.

865

866

# References

Artuso, I., Lidonnici, M.R., Altamura, S., Mandelli, G., Pettinato, M., Muckenthaler, M.U., Silvestri, L., Ferrari, G., Camaschella, C., and Nai, A. (2018). Transferrin Receptor 2 is a potential novel therapeutic target for beta-thalassemia: evidence from a murine model. Blood.

Astle, W.J., Elding, H., Jiang, T., Allen, D., Ruklisa, D., Mann, A.L., Mead, D., Bouman, H., Riveros-Mckay, F., Kostadima, M.A., *et al.* (2016). The Allelic Landscape of Human Blood Cell Trait Variation and Links to Common Complex Disease. Cell *167*, 1415-1429.e1419.

Basak, A., Hancarova, M., Ulirsch, J.C., Balci, T.B., Trkova, M., Pelisek, M., Vlckova, M., Muzikova, K., Cermak, J., Trka, J., *et al.* (2015). BCL11A deletions result in fetal hemoglobin persistence and neurodevelopmental alterations. The Journal of clinical investigation *125*, 2363-2368.

Bennett, V., and Stenbuck, P.J. (1979). Identification and partial purification of ankyrin, the high affinity membrane attachment site for human erythrocyte spectrin. J Biol Chem *254*, 2533-2541.

Boyle, E.A., Li, Y.I., and Pritchard, J.K. (2017). An Expanded View of Complex Traits: From Polygenic to Omnigenic. Cell *169*, 1177-1186.

Claussnitzer, M., Dankel, S.N., Kim, K.H., Quon, G., Meuleman, W., Haugen, C., Glunk, V., Sousa, I.S., Beaudry, J.L., Puviindran, V., *et al.* (2015). FTO Obesity Variant Circuitry and Adipocyte Browning in Humans. The New England journal of medicine *373*, 895-907.

Corces, M.R., Buenrostro, J.D., Wu, B., Greenside, P.G., Chan, S.M., Koenig, J.L., Snyder, M.P., Pritchard, J.K., Kundaje, A., Greenleaf, W.J., *et al.* (2016). Lineage-specific and single-cell chromatin accessibility charts human hematopoiesis and leukemia evolution. Nat Genet *48*, 1193-1203.

Eden, E., Lipson, D., Yogev, S., and Yakhini, Z. (2007). Discovering motifs in ranked lists of DNA sequences. PLoS Comput Biol *3*, e39.

Eden, E., Navon, R., Steinfeld, I., Lipson, D., and Yakhini, Z. (2009). GOrilla: a tool for discovery and visualization of enriched GO terms in ranked gene lists. BMC Bioinformatics *10*, 48.

Forejtnikova, H., Vieillevoye, M., Zermati, Y., Lambert, M., Pellegrino, R.M., Guihard, S., Gaudry, M., Camaschella, C., Lacombe, C., Roetto, A., *et al.* (2010). Transferrin receptor 2 is a component of the erythropoietin receptor complex and is required for efficient erythropoiesis. Blood *116*, 5357-5367.

Fulco, C.P., Munschauer, M., Anyoha, R., Munson, G., Grossman, S.R., Perez, E.M., Kane, M., Cleary, B., Lander, E.S., and Engreitz, J.M. (2016). Systematic mapping of functional enhancer-promoter connections with CRISPR interference. Science *354*, 769-773.

Galarneau, G., Palmer, C.D., Sankaran, V.G., Orkin, S.H., Hirschhorn, J.N., and Lettre, G. (2010). Fine-mapping at three loci known to affect fetal hemoglobin levels explains additional genetic variation. Nat Genet *42*, 1049-1051.

903  Gazda, H.T., Preti, M., Sheen, M.R., O'Donohue, M.F., Vlachos, A., Davies, S.M., Kattamis, A.,
904  Doherty, L., Landowski, M., Buros, C*., et al.* (2012). Frameshift mutation in p53 regulator RPL26
905  is associated with multiple physical abnormalities and a specific pre-rRNA processing defect in
906  Diamond-Blackfan anemia. Hum Mutat *33*, 1037-1044.

907  Giani, F.C., Fiorini, C., Wakabayashi, A., Ludwig, L.S., Salem, R.M., Jobaliya, C.D., Regan,
908  S.N., Ulirsch, J.C., Liang, G., Steinberg-Shemer, O*., et al.* (2016). Targeted Application of
909  Human Genetic Variation Can Improve Red Blood Cell Production from Stem Cells. Cell Stem
910  Cell *18*, 73-78.

911  Glick, B.S., and Rothman, J.E. (1987). Possible role for fatty acyl-coenzyme A in intracellular
912  protein transport. Nature *326*, 309-312.

913  Goardon, N., Lambert, J.A., Rodriguez, P., Nissaire, P., Herblot, S., Thibault, P., Dumenil, D.,
914  Strouboulis, J., Romeo, P.H., and Hoang, T. (2006). ETO2 coordinates cellular proliferation and
915  differentiation during erythropoiesis. Embo j *25*, 357-366.

916  Gozani, O., Feld, R., and Reed, R. (1996). Evidence that sequence-independent binding of
917  highly conserved U2 snRNP proteins upstream of the branch site is required for assembly of
918  spliceosomal complex A. Genes Dev *10*, 233-243.

919  Gozani, O., Potashkin, J., and Reed, R. (1998). A potential role for U2AF-SAP 155 interactions
920  in recruiting U2 snRNP to the branch site. Molecular and cellular biology *18*, 4752-4760.

921  Guo, M.H., Nandakumar, S.K., Ulirsch, J.C., Zekavat, S.M., Buenrostro, J.D., Natarajan, P.,
922  Salem, R.M., Chiarle, R., Mitt, M., and Kals, M. (2016). Comprehensive population-based
923  genome sequencing provides insight into hematopoietic regulatory mechanisms. Proceedings of
924  the National Academy of Sciences, 201619052.

925  Gupta, R.M., Hadaya, J., Trehan, A., Zekavat, S.M., Roselli, C., Klarin, D., Emdin, C.A.,
926  Hilvering, C.R.E., Bianchi, V., Mueller, C*., et al.* (2017). A Genetic Variant Associated with Five
927  Vascular Diseases Is a Distal Regulator of Endothelin-1 Gene Expression. Cell *170*, 522-
928  533.e515.

929  Hu, J., Liu, J., Xue, F., Halverson, G., Reid, M., Guo, A., Chen, L., Raza, A., Galili, N., Jaffray, J*.,*
930  *et al.* (2013). Isolation and functional characterization of human erythroblasts at distinct stages:
931  implications for understanding of normal and disordered erythropoiesis in vivo. Blood *121*, 3246-
932  3253.

933  Huang, H., Fang, M., Jostins, L., Umicevic Mirkov, M., Boucher, G., Anderson, C.A., Andersen,
934  V., Cleynen, I., Cortes, A., Crins, F*., et al.* (2017). Fine-mapping inflammatory bowel disease loci
935  to single-variant resolution. Nature *547*, 173-178.

936  Jing, H., Vakoc, C.R., Ying, L., Mandat, S., Wang, H., Zheng, X., and Blobel, G.A. (2008).
937  Exchange of GATA factors mediates transitions in looped chromatin organization at a
938  developmentally regulated gene locus. Mol Cell *29*, 232-242.

939  Katz, Y., Wang, E.T., Airoldi, E.M., and Burge, C.B. (2010). Analysis and design of RNA
940  sequencing experiments for identifying isoform regulation. Nat Methods *7*, 1009-1015.

941  Khajuria, R.K., Munschauer, M., Ulirsch, J.C., Fiorini, C., Ludwig, L.S., McFarland, S.K.,
942  Abdulhay, N.J., Specht, H., Keshishian, H., Mani, D.R., *et al.* (2018). Ribosome Levels
943  Selectively Regulate Translation and Lineage Commitment in Human Hematopoiesis. Cell *173*,
944  90-103.e119.

945  Lareau, C.A., Ulirsch, J.C., Bao, E.L., Ludwig, L.S., Guo, M.H., Benner, C., Satpathy, A.T.,
946  Salem, R., Hirschhorn, J.N., Finucane, H.K., *et al.* (2018). Interrogation of human hematopoiesis
947  at single-cell and single-variant resolution.

948  Li, W., Koster, J., Xu, H., Chen, C.H., Xiao, T., Liu, J.S., Brown, M., and Liu, X.S. (2015). Quality
949  control, modeling, and visualization of CRISPR screens with MAGeCK-VISPR. Genome biology
950  *16*, 281.

951  Liu, N., Hargreaves, V.V., Zhu, Q., Kurland, J.V., Hong, J., Kim, W., Sher, F., Macias-Trevino, C.,
952  Rogers, J.M., Kurita, R., *et al.* (2018). Direct Promoter Repression by BCL11A Controls the
953  Fetal to Adult Hemoglobin Switch. Cell *173*, 430-442.e417.

954  Love, M.I., Huber, W., and Anders, S. (2014). Moderated estimation of fold change and
955  dispersion for RNA-seq data with DESeq2. Genome Biol *15*, 550.

956  Ludwig, L.S., Gazda, H.T., Eng, J.C., Eichhorn, S.W., Thiru, P., Ghazvinian, R., George, T.I.,
957  Gotlib, J.R., Beggs, A.H., Sieff, C.A., *et al.* (2014). Altered translation of GATA1 in Diamond-
958  Blackfan anemia. Nature medicine *20*, 748-753.

959  Mandegar, M.A., Huebsch, N., Frolov, E.B., Shin, E., Truong, A., Olvera, M.P., Chan, A.H.,
960  Miyaoka, Y., Holmes, K., Spencer, C.I., *et al.* (2016). CRISPR Interference Efficiently Induces
961  Specific and Reversible Gene Silencing in Human iPSCs. Cell Stem Cell *18*, 541-553.

962  McIver, S.C., Kang, Y.A., DeVilbiss, A.W., O'Driscoll, C.A., Ouellette, J.N., Pope, N.J.,
963  Camprecios, G., Chang, C.J., Yang, D., Bouhassira, E.E., *et al.* (2014). The exosome complex
964  establishes a barricade to erythroid maturation. Blood *124*, 2285-2297.

965  Moffat, J., Grueneberg, D.A., Yang, X., Kim, S.Y., Kloepfer, A.M., Hinkle, G., Piqani, B.,
966  Eisenhaure, T.M., Luo, B., Grenier, J.K., *et al.* (2006). A lentiviral RNAi library for human and
967  mouse genes applied to an arrayed viral high-content screen. Cell *124*, 1283-1298.

968  Mohanan, V., Nakata, T., Desch, A.N., Levesque, C., Boroughs, A., Guzman, G., Cao, Z.,
969  Creasey, E., Yao, J., Boucher, G., *et al.* (2018). C1orf106 is a colitis risk gene that regulates
970  stability of epithelial adherens junctions. Science *359*, 1161-1166.

971  Moniz, H., Gastou, M., Leblanc, T., Hurtaud, C., Cretien, A., Lecluse, Y., Raslova, H., Larghero,
972  J., Croisille, L., Faubladier, M., *et al.* (2012). Primary hematopoietic cells from DBA patients with
973  mutations in RPL11 and RPS19 genes exhibit distinct erythroid phenotype in vitro. Cell Death
974  Dis *3*, e356.

975  Montgomery, S.B., and Dermitzakis, E.T. (2011). From expression QTLs to personalized
976  transcriptomics. Nature reviews Genetics *12*, 277-282.

977    Mucenski, M.L., McLain, K., Kier, A.B., Swerdlow, S.H., Schreiner, C.M., Miller, T.A., Pietryga,
978    D.W., Scott, W.J., Jr., and Potter, S.S. (1991). A functional c-myb gene is required for normal
979    murine fetal hepatic hematopoiesis. Cell *65*, 677-689.

980    Musunuru, K., Strong, A., Frank-Kamenetsky, M., Lee, N.E., Ahfeldt, T., Sachs, K.V., Li, X., Li,
981    H., Kuperwasser, N., Ruda, V.M.*, et al.* (2010). From noncoding variant to phenotype via SORT1
982    at the 1p13 cholesterol locus. Nature *466*, 714-719.

983    Nai, A., Lidonnici, M.R., Rausa, M., Mandelli, G., Pagani, A., Silvestri, L., Ferrari, G., and
984    Camaschella, C. (2015). The second transferrin receptor regulates red blood cell production in
985    mice. Blood *125*, 1170-1179.

986    Niemi, M.E.K., Martin, H.C., Rice, D.L., Gallone, G., Gordon, S., Kelemen, M., McAloney, K.,
987    McRae, J., Radford, E.J., Yu, S.*, et al.* (2018). Common genetic variants contribute to risk of
988    rare severe neurodevelopmental disorders. Nature *562*, 268-271.

989    Obeng, E.A., Chappell, R.J., Seiler, M., Chen, M.C., Campagna, D.R., Schmidt, P.J., Schneider,
990    R.K., Lord, A.M., Wang, L., Gambe, R.G.*, et al.* (2016). Physiologic Expression of Sf3b1(K700E)
991    Causes Impaired Erythropoiesis, Aberrant Splicing, and Sensitivity to Therapeutic Spliceosome
992    Modulation. Cancer cell *30*, 404-417.

993    Papaemmanuil, E., Cazzola, M., Boultwood, J., Malcovati, L., Vyas, P., Bowen, D., Pellagatti, A.,
994    Wainscoat, J.S., Hellstrom-Lindberg, E., Gambacorti-Passerini, C.*, et al.* (2011). Somatic SF3B1
995    mutation in myelodysplasia with ring sideroblasts. The New England journal of medicine *365*,
996    1384-1395.

997    Patro, R., Duggal, G., Love, M.I., Irizarry, R.A., and Kingsford, C. (2017). Salmon provides fast
998    and bias-aware quantification of transcript expression. Nat Methods *14*, 417-419.

999    Peters, L.L., Shivdasani, R.A., Liu, S.C., Hanspal, M., John, K.M., Gonzalez, J.M., Brugnara, C.,
1000   Gwynn, B., Mohandas, N., Alper, S.L.*, et al.* (1996). Anion exchanger 1 (band 3) is required to
1001   prevent erythrocyte membrane surface loss but not to form the membrane skeleton. Cell *86*,
1002   917-927.

1003   Pimentel, H., Parra, M., Gee, S.L., Mohandas, N., Pachter, L., and Conboy, J.G. (2016). A
1004   dynamic intron retention program enriched in RNA processing genes regulates gene expression
1005   during terminal erythropoiesis. Nucleic Acids Res *44*, 838-851.

1006   Riba, A., Emmenlauer, M., Chen, A., Sigoillot, F., Cong, F., Dehio, C., Jenkins, J., and Zavolan,
1007   M. (2017). Explicit Modeling of siRNA-Dependent On- and Off-Target Repression Improves the
1008   Interpretation of Screening Results. Cell systems *4*, 182-193.e184.

1009   Rossi, A., Kontarakis, Z., Gerri, C., Nolte, H., Holper, S., Kruger, M., and Stainier, D.Y. (2015).
1010   Genetic compensation induced by deleterious mutations but not gene knockdowns. Nature *524*,
1011   230-233.

1012   Rossin, E.J., Lage, K., Raychaudhuri, S., Xavier, R.J., Tatar, D., Benita, Y., Cotsapas, C., and
1013   Daly, M.J. (2011). Proteins encoded in genomic regions associated with immune-mediated
1014   disease physically interact and suggest underlying biology. PLoS genetics *7*, e1001273.

1015   Rund, D., and Rachmilewitz, E. (2005). Beta-thalassemia. N Engl J Med *353*, 1135-1146.

1016   Sankaran, V.G., Joshi, M., Agrawal, A., Schmitz-Abe, K., Towne, M.C., Marinakis, N., Markianos,
1017   K., Berry, G.T., and Agrawal, P.B. (2013). Rare complete loss of function provides insight into a
1018   pleiotropic genome-wide association study locus. Blood *122*, 3845-3847.

1019   Sankaran, V.G., Ludwig, L.S., Sicinska, E., Xu, J., Bauer, D.E., Eng, J.C., Patterson, H.C.,
1020   Metcalf, R.A., Natkunam, Y., Orkin, S.H.*, et al.* (2012). Cyclin D3 coordinates the cell cycle
1021   during differentiation to regulate erythrocyte size and number. Genes Dev *26*, 2075-2087.

1022   Sankaran, V.G., Menne, T.F., Scepanovic, D., Vergilio, J.A., Ji, P., Kim, J., Thiru, P., Orkin, S.H.,
1023   Lander, E.S., and Lodish, H.F. (2011). MicroRNA-15a and -16-1 act via MYB to elevate fetal
1024   hemoglobin expression in human trisomy 13. Proceedings of the National Academy of Sciences
1025   of the United States of America *108*, 1519-1524.

1026   Sankaran, V.G., Menne, T.F., Xu, J., Akie, T.E., Lettre, G., Van Handel, B., Mikkola, H.K.,
1027   Hirschhorn, J.N., Cantor, A.B., and Orkin, S.H. (2008). Human fetal hemoglobin expression is
1028   regulated by the developmental stage-specific repressor BCL11A. Science *322*, 1839-1842.

1029   Simeonov, D.R., Gowen, B.G., Boontanrart, M., Roth, T.L., Gagnon, J.D., Mumbach, M.R.,
1030   Satpathy, A.T., Lee, Y., Bray, N.L., Chan, A.Y.*, et al.* (2017). Discovery of stimulation-responsive
1031   immune enhancers with CRISPR activation. Nature *549*, 111-115.

1032   Smemo, S., Tena, J.J., Kim, K.H., Gamazon, E.R., Sakabe, N.J., Gomez-Marin, C., Aneas, I.,
1033   Credidio, F.L., Sobreira, D.R., Wasserman, N.F.*, et al.* (2014). Obesity-associated variants within
1034   FTO form long-range functional connections with IRX3. Nature *507*, 371-375.

1035   Szklarczyk, D., Morris, J.H., Cook, H., Kuhn, M., Wyder, S., Simonovic, M., Santos, A.,
1036   Doncheva, N.T., Roth, A., Bork, P.*, et al.* (2017). The STRING database in 2017: quality-
1037   controlled protein-protein association networks, made broadly accessible. Nucleic Acids Res *45*,
1038   D362-d368.

1039   Tewhey, R., Kotliar, D., Park, D.S., Liu, B., Winnicki, S., Reilly, S.K., Andersen, K.G., Mikkelsen,
1040   T.S., Lander, E.S., and Schaffner, S.F. (2016). Direct identification of hundreds of expression-
1041   modulating variants using a multiplexed reporter assay. Cell *165*, 1519-1529.

1042   Thomsen, S.K., Ceroni, A., van de Bunt, M., Burrows, C., Barrett, A., Scharfmann, R., Ebner, D.,
1043   McCarthy, M.I., and Gloyn, A.L. (2016). Systematic Functional Characterization of Candidate
1044   Causal Genes for Type 2 Diabetes Risk Variants. Diabetes *65*, 3805-3811.

1045   Ting, P.Y., Parker, A.E., Lee, J.S., Trussell, C., Sharif, O., Luna, F., Federe, G., Barnes, S.W.,
1046   Walker, J.R., Vance, J.*, et al.* (2018). Guide Swap enables genome-scale pooled CRISPR-Cas9
1047   screening in human primary cells. Nat Methods *15*, 941-946.

1048   Tsherniak, A., Vazquez, F., Montgomery, P.G., Weir, B.A., Kryukov, G., Cowley, G.S., Gill, S.,
1049   Harrington, W.F., Pantel, S., Krill-Burger, J.M.*, et al.* (2017). Defining a Cancer Dependency
1050   Map. Cell *170*, 564-576.e516.

1051    Ulirsch, J.C., Nandakumar, S.K., Wang, L., Giani, F.C., Zhang, X., Rogov, P., Melnikov, A.,
1052    McDonel, P., Do, R., Mikkelsen, T.S.*, et al.* (2016). Systematic Functional Dissection of Common
1053    Genetic Variation Affecting Red Blood Cell Traits. Cell *165*, 1530-1545.

1054    Ulirsch, J.C., Verboon, J.M., Kazerounian, S., Guo, M.H., Yuan, D., Ludwig, L.S., Handsaker,
1055    R.E., Abdulhay, N.J., Fiorini, C., Genovese, G.*, et al.* (2018). The Genetic Landscape of
1056    Diamond-Blackfan Anemia.

1057    van der Harst, P., Zhang, W., Mateo Leach, I., Rendon, A., Verweij, N., Sehmi, J., Paul, D.S.,
1058    Elling, U., Allayee, H., Li, X.*, et al.* (2012). Seventy-five genetic loci influencing the human red
1059    blood cell. Nature *492*, 369-375.

1060    Veyrieras, J.B., Kudaravalli, S., Kim, S.Y., Dermitzakis, E.T., Gilad, Y., Stephens, M., and
1061    Pritchard, J.K. (2008). High-resolution mapping of expression-QTLs yields insight into human
1062    gene regulation. PLoS genetics *4*, e1000214.

1063    Vockley, C.M., Guo, C., Majoros, W.H., Nodzenski, M., Scholtens, D.M., Hayes, M.G., Lowe,
1064    W.L., Jr., and Reddy, T.E. (2015). Massively parallel quantification of the regulatory effects of
1065    noncoding genetic variation in a human cohort. Genome research *25*, 1206-1214.

1066    Wang, T., Birsoy, K., Hughes, N.W., Krupczak, K.M., Post, Y., Wei, J.J., Lander, E.S., and
1067    Sabatini, D.M. (2015). Identification and characterization of essential genes in the human
1068    genome. Science *350*, 1096-1101.
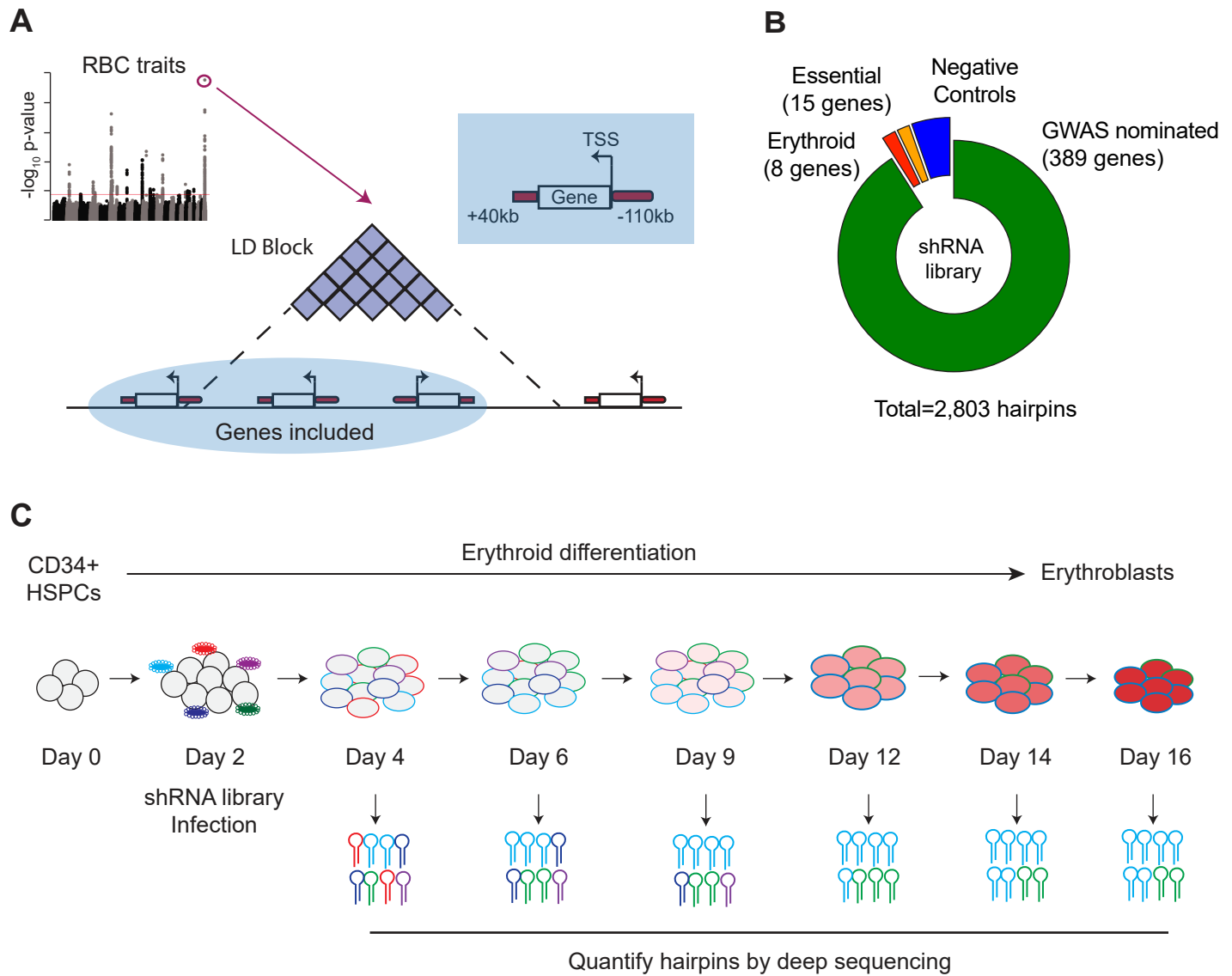
1069    Wang, X., Angelis, N., and Thein, S.L. (2018). MYB - A regulatory factor in hematopoiesis. Gene
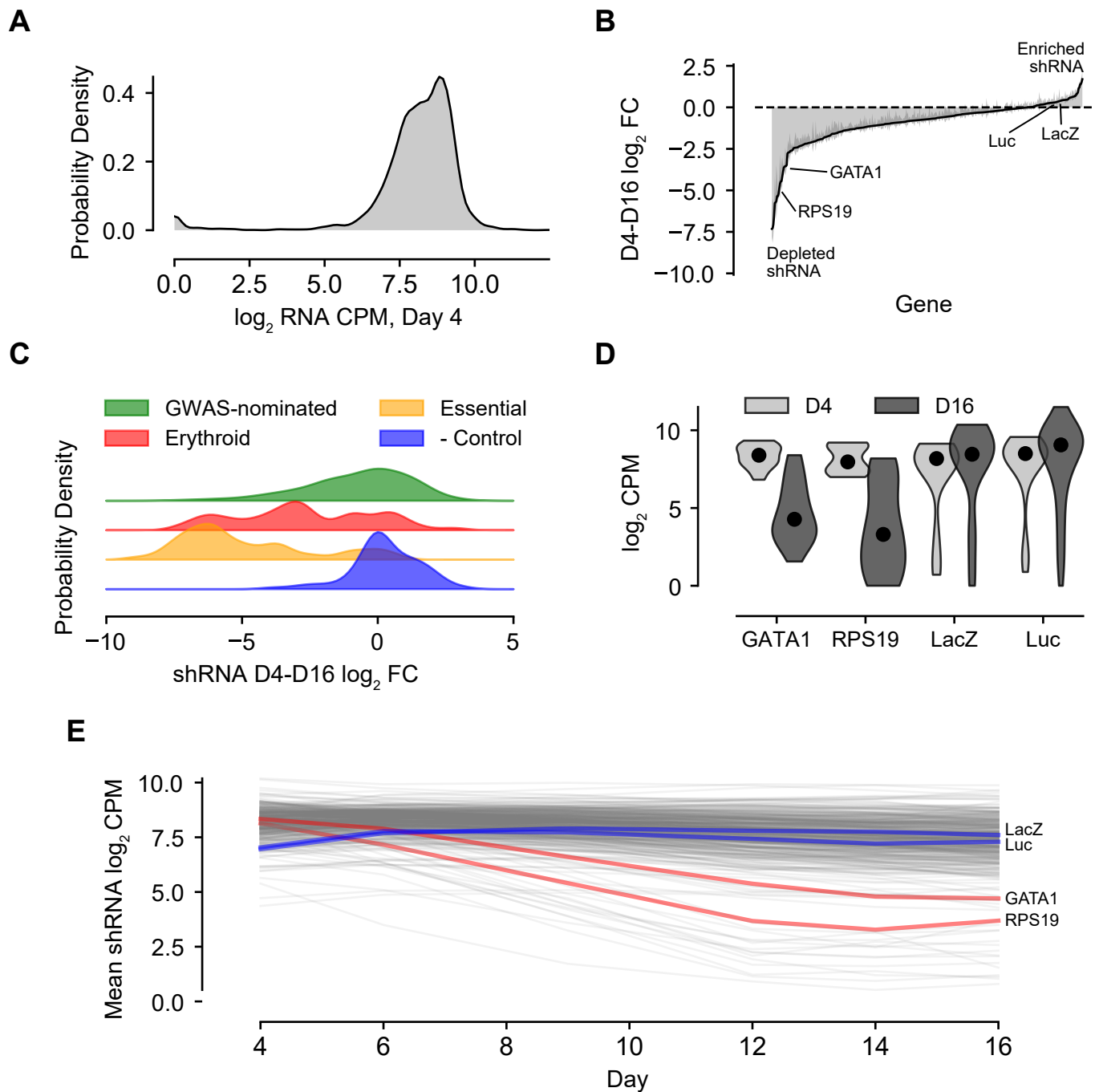1070    *665*, 6-17.

1071    Whalen, S., and Pollard, K.S. (2018). Most regulatory interactions are not in linkage
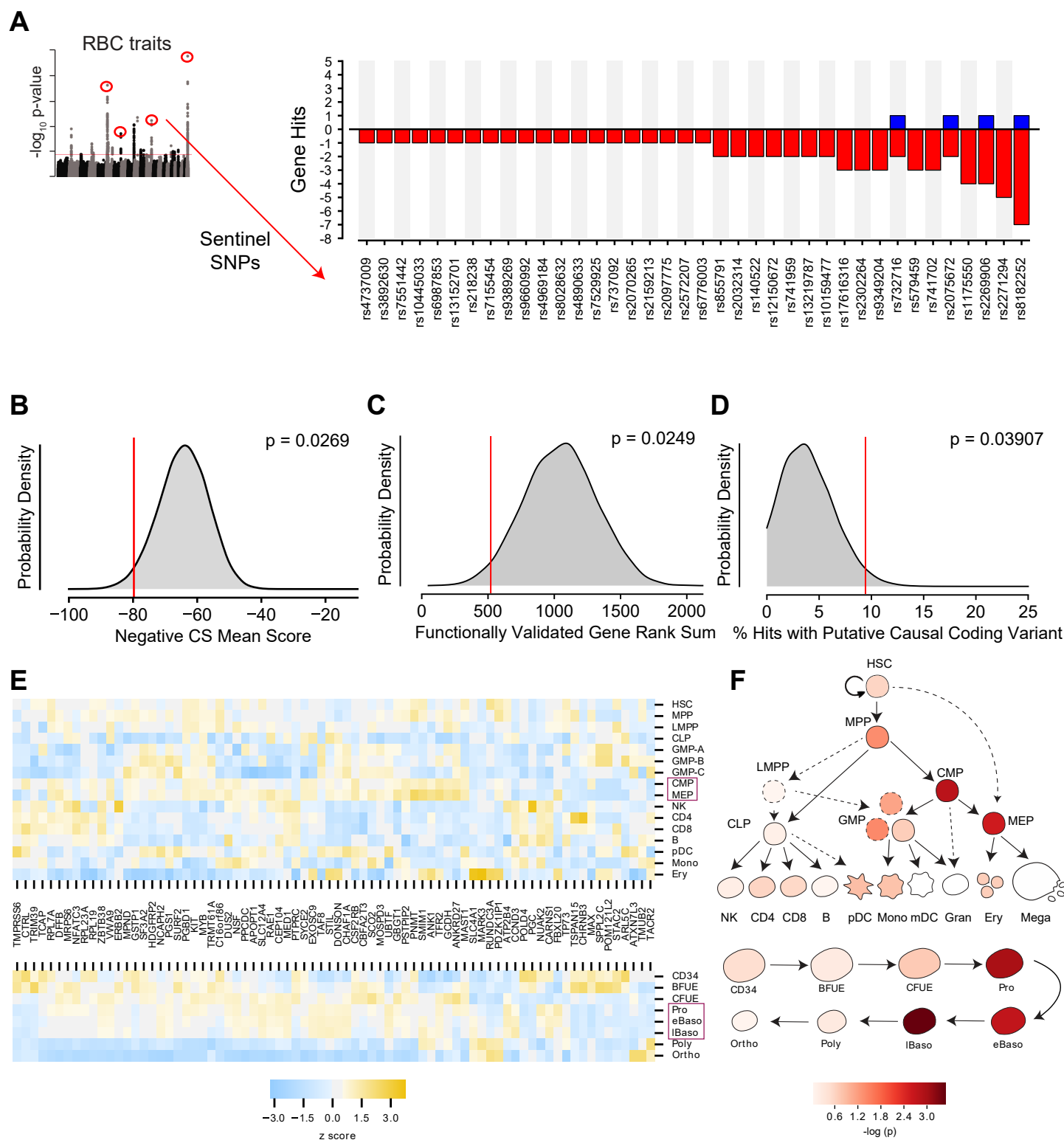1072    disequilibrium (Cold Spring Harbor Laboratory).

1073    Willer, C.J., Schmidt, E.M., Sengupta, S., Peloso, G.M., Gustafsson, S., Kanoni, S., Ganna, A.,
1074    Chen, J., Buchkovich, M.L., Mora, S.*, et al.* (2013). Discovery and refinement of loci associated
1075    with lipid levels. Nat Genet *45*, 1274-1283.

1076    Yan, H., Hale, J., Jaffray, J., Li, J., Wang, Y., Huang, Y., An, X., Hillyer, C., Wang, N., Kinet, S.*,
1077    et al.* (2018). Developmental differences between neonatal and adult human erythropoiesis. Am
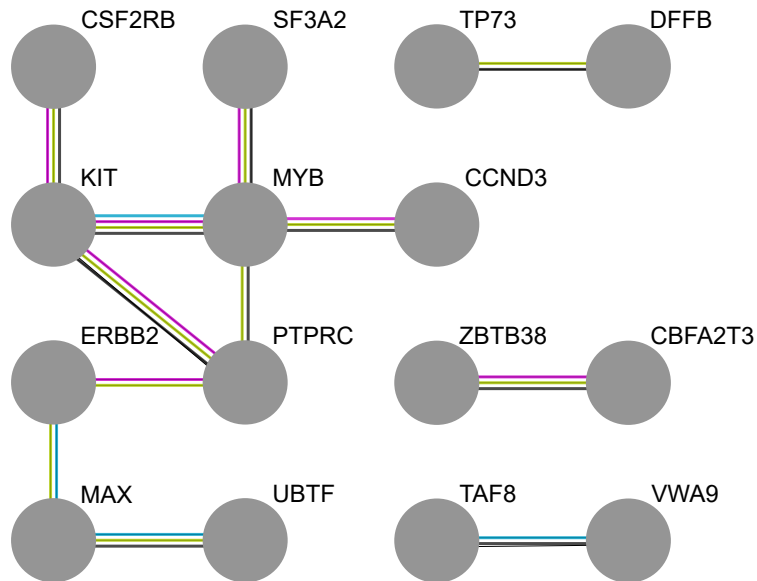1078    J Hematol *93*, 494-503.
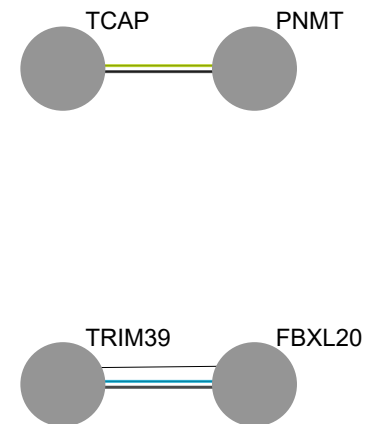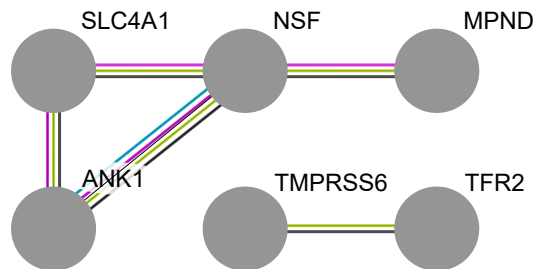1079

Figure 1



**A**

RBC traits

-log₁₀ p-value

LD Block

Genes included

TSS

Gene

+40kb          -110kb

**B**

Essential
(15 genes)

Negative
Controls

Erythroid
(8 genes)

GWAS nominated
(389 genes)

shRNA
library

Total=2,803 hairpins

**C**

CD34+
HSPCs

Erythroid differentiation

Erythroblasts

Day 0    Day 2    Day 4    Day 6    Day 9    Day 12    Day 14    Day 16

shRNA library
Infection

Quantify hairpins by deep sequencing

Figure 2

Figure 3

Figure 4

Figure 5

Figure 6

Figure S1

Figure S2