

# Systematic documentation and analysis of human genetic variation in hemoglobinopathies using the microattribution approach

Belinda Giardine<sup>1,37\*</sup>, Joseph Borg<sup>2-4,37</sup>, Douglas R Higgs<sup>5</sup>, Kenneth R Peterson<sup>6</sup>, Sjaak Philipsen<sup>7</sup>, Donna Maglott<sup>8</sup>, Belinda K Singleton<sup>9</sup>, David J Anstee<sup>9</sup>, A Nazli Basak<sup>10</sup>, Barnaby Clark<sup>11</sup>, Flavia C Costa<sup>6</sup>, Paula Faustino<sup>12</sup>, Halyna Fedosyuk<sup>6</sup>, Alex E Felice<sup>3,4</sup>, Alain Francina<sup>13</sup>, Renzo Galanello<sup>14</sup>, Monica V E Gallivan<sup>15</sup>, Marianthi Georgitsi<sup>16</sup>, Richard J Gibbons<sup>5</sup>, Piero C Giordano<sup>17</sup>, Cornelis L Harteveld<sup>17</sup>, James D Hoyer<sup>18</sup>, Martin Jarvis<sup>19</sup>, Philippe Joly<sup>13</sup>, Emmanuel Kanavakis<sup>20</sup>, Panagoula Kollia<sup>21</sup>, Stephan Menzel<sup>11</sup>, Webb Miller<sup>1</sup>, Kamran Moradkhani<sup>22</sup>, John Old<sup>23</sup>, Adamantia Papachatzopoulou<sup>24</sup>, Manoussos N Papadakis<sup>25</sup>, Petros Papadopoulos<sup>7</sup>, Sonja Pavlovic<sup>26</sup>, Lucia Perseu<sup>27</sup>, Milena Radmilovic<sup>26</sup>, Cathy Riemer<sup>1</sup>, Stefania Satta<sup>14</sup>, Iris Schrijver<sup>28</sup>, Maja Stojiljkovic<sup>26</sup>, Swee Lay Thein<sup>11</sup>, Jan Traeger-Synodinos<sup>20</sup>, Ray Tully<sup>8</sup>, Takahito Wada<sup>29</sup>, John S Wayne<sup>30,31</sup>, Claudia Wiemann<sup>32</sup>, Branka Zukic<sup>26</sup>, David H K Chui<sup>33,34</sup>, Henri Wajcman<sup>22,35</sup>, Ross C Hardison<sup>1,36</sup> & George P Patrinos<sup>16</sup>

**We developed a series of interrelated locus-specific databases to store all published and unpublished genetic variation related to hemoglobinopathies and thalassemia and implemented microattribution to encourage submission of unpublished observations of genetic variation to these public repositories. A total of 1,941 unique genetic variants in 37 genes, encoding globins and other erythroid proteins, are currently documented in these databases, with reciprocal attribution of microcitations to data contributors. Our project provides the first example of implementing microattribution to incentivise submission of all known genetic variation in a defined system. It has demonstrably increased the reporting of human variants, leading to a comprehensive online resource for systematically describing human genetic variation in the globin genes and other genes contributing to hemoglobinopathies and thalassemias. The principles established here will serve as a model for other systems and for the analysis of other common and/or complex human genetic diseases.**

Since completion of the human genome project, a major aim in the field of genetics has been to determine how individual genomes differ from each other and how these differences explain variation in phenotype. However, it often remains unclear which variants cause changes in phenotype and which are phenotype neutral; furthermore, in many instances, the mechanisms by which variants cause changes in gene expression and phenotypes remain unknown. To address this, DNA sequence data will need to be matched

with well-defined phenotypes to make meaningful connections between structure, function and mechanism.

A potential hurdle to this approach is how to encourage 'phenotypers' to report their observations. After the initial excitement during the 1980s and 1990s of identifying disease-causing molecular defects and the mechanisms by which they arise, enthusiasm in this area has declined such that it has become increasingly difficult to report small numbers of human variants in scientific journals. Consequently, many new variants associated with well-defined phenotypes and, equally important, variants which cause no change in phenotype remain unreported. Inevitably, a large amount of potentially valuable information remains inaccessible.

To overcome this problem, we implemented a process for capturing such information with the incentive of microattribution, whereby the contribution of those individuals collecting new detailed genotype and phenotype data is positively encouraged and appropriately acknowledged<sup>1</sup>. We have applied the microattribution approach to inherited disorders affecting either the structure of hemoglobin (such as sickle cell disease (SCD)) or the levels and balance of globin chain production (the thalassemias). We also included variants that cause hereditary persistence of fetal hemoglobin (HPFH), a condition associated with increased production of  $\gamma$ -globin which ameliorates the clinical endpoints of SCD and  $\beta$ -thalassemia. The hemoglobinopathies and thalassemias are among the commonest inherited disorders in humans. Variants of the globin-encoding genes, residing in the  $\alpha$ -like and  $\beta$ -like globin gene clusters, have provided key insights into the principles underlying human molecular genetics since the discipline was established in the 1950s (ref. 2).

\*A full list of author affiliations appear at the end of the paper.

Although most hemoglobinopathies are classic monogenic disorders affecting structural genes, globin gene expression is the end product of a complex regulatory network (transcriptional and epigenetic) that emerges during terminal erythroid differentiation. Consequently, globin gene expression may also be affected by *trans*-acting mutations. Examples of such mutations were initially found in families with rare syndromal disorders, of which  $\alpha$ -thalassemia was one component (for example, ATR-X (MIM301040) and ATMDS (MIM300448) syndromes)<sup>3,4</sup>. Similarly, trichothiodystrophy (MIM300448) was shown to be associated with  $\beta$ -thalassemia due to mutations in the XPD component of the general transcription factor complex TFIID<sup>5</sup>. The association of X-linked thrombocytopenia with  $\beta$ -thalassemia identified a mutation of the erythroid-specific transcription factor GATA-1 (ref. 6), and recently, systematic analysis of subjects with unexplained HPFH has identified mutations in the KLF1 erythroid transcription factor<sup>7</sup>. Finally, the implementation of genome-wide association studies searching for quantitative trait loci that influence the level of fetal hemoglobin (HbF) has revealed several important regulators of *HBG1* and *HBG2* gene expression, including the *HBSIL-MYB*<sup>8</sup> and *BCL11A* loci<sup>9,10</sup> on chromosomes 6 and 2, respectively. As genetic variations in the genes within the erythroid network are investigated in further detail, we anticipate many more discoveries of *trans*-acting mutations that may provide target pathways for manipulating globin gene expression to ameliorate the symptoms of thalassemia and SCD. Therefore it is important that an effective database be created to accommodate all of the mutations affecting the globin genes and the network regulating their expression.

Here we report the first example of implementing microattribution to systematically document genetic variation leading to human genetic disorders, using hemoglobinopathies and thalassemias as an example. Furthermore, we demonstrate that microattribution can incentivise data contribution and, importantly, show how an integrated human variant database (including the recently acquired microattribution data) can provide key insights into human genetic diseases. Microattribution provides an important mechanism and incentive for researchers to report all variants within a specific gene or disease network. Following the principles established for the globin disorders, these databases should provide a key resource for understanding the molecular pathology of human genetic diseases.

### Developing the microattribution process

To ensure that all natural mutations and their associated phenotypes are accurately and efficiently recorded, we comprehensively documented genotype and phenotype information in individuals with globin disorders in a series of interrelated locus-specific databases (LSDBs). Traditionally, credit has been given to discoverers of genetic variants through citations of their publications describing the variants. However, the increased rate of discovery through re-sequencing efforts far exceeds the capacity of citations of individual publications to give adequate credit. In order to be used effectively by the community, published variants are deposited into databases such as those described here; nevertheless, many variants may still not be published. Alternatively, variants may be discovered in large-scale collaborative projects. Credit can be given to the discoverers of the variants deposited in databases through the new process of microattribution<sup>1</sup>. Each variant used in a paper is listed in four microattribution tables with its accession number and with unique IDs for the discoverers, or 'authors', of the variant. In this paper, we have applied 'microcitations' to hemoglobinopathy-associated variants in order to provide incentives to data producers to deposit all of their data in these public resources<sup>1</sup>. Depositing the microattribution tables in a central repository

(for example, NCBI) provides a venue for quantitative microcitations for every unique author. Using this approach (first implemented in 2010), there has been a marked increase in the number of reported variants in the globin gene network (**Supplementary Fig. 1**).

### Implementing microattribution

All genetic variation data have been collected and documented in the HbVar database of hemoglobin variants and thalassemia mutations<sup>11</sup> and the Leiden Open-Access Variation Database (LOVD)-based LSDBs for the other erythroid proteins<sup>12</sup> (**Supplementary Note**) with appropriate attribution of the data contributors. These variants are reported in publicly available microattribution tables (also provided in **Supplementary Table 1**) that have been centrally deposited in NCBI (**Supplementary Fig. 2**). Each microattribution table has different information related to submission to the central depository, microattribution, phenotype and allele frequency (**Supplementary Note**).

In this protocol, data submitters directly contribute variants leading to hemoglobinopathies to HbVar and in return obtain direct microattribution credit. These variants have been recorded with researcher IDs and in the case of previously published variants, the corresponding PubMed ID was also used (**Supplementary Fig. 2**). To date, 232 variants have been directly submitted to HbVar without being published in a peer-reviewed journal, some of which have been deposited with more than one researcher ID. Seventy-six variants were 'orphan', that is, variants for which there was neither a PubMed ID nor a researcher ID, all of which were variants initially deposited to HbVar in the year 2000 and for which either valid contact details for the variant contributors was lacking or the contributor(s) failed to respond to our invitation. These variants have been deposited with an HbVar researcher ID.

For all unpublished variants directly contributed to HbVar by the microattribution process, a very stringent evaluation of the information submitted takes place. Contributed variant data are evaluated by curators, all of whom are senior scientists with extensive editorial experience, especially in the field of hemoglobinopathies. The curators directly contact the data contributors, if needed, for clarifications related to issues pertaining to phenotypic description, method of variant identification, ethnicity of the individual with the variant, allele frequency and so on. Upon acceptance, contributed data become part of the main HbVar data collection recorded with the contributor(s) researcher ID.

Although microattribution can operate locally (within journals and databases each reporting quantitative citation of accessions), depositing the microattribution tables in a central repository of cited accessions (for example, NCBI or European Bioinformatics Institute (EBI)) allows the central registry to be mined for citations associated with unique author identities and with each author's publications and database entries. For the purpose of our project, we have chosen to deposit the microattribution tables in NCBI, and a copy of these tables is also deposited in Nature Publishing Group's central database.

### Mining the databases

In the case of globin gene disorders, many variants were conventionally reported in genetics journals, and these variants identified and/or elucidated many mechanisms underlying key aspects of gene regulation in *cis* (for example, promoters, enhancers, silencers, mRNA processing signals and translational signals) and in *trans* (for example, transcription factors, chromatin remodeling factors and protein chaperones)<sup>2</sup>. Furthermore, these variants helped to establish the molecular mechanisms underlying human genetic

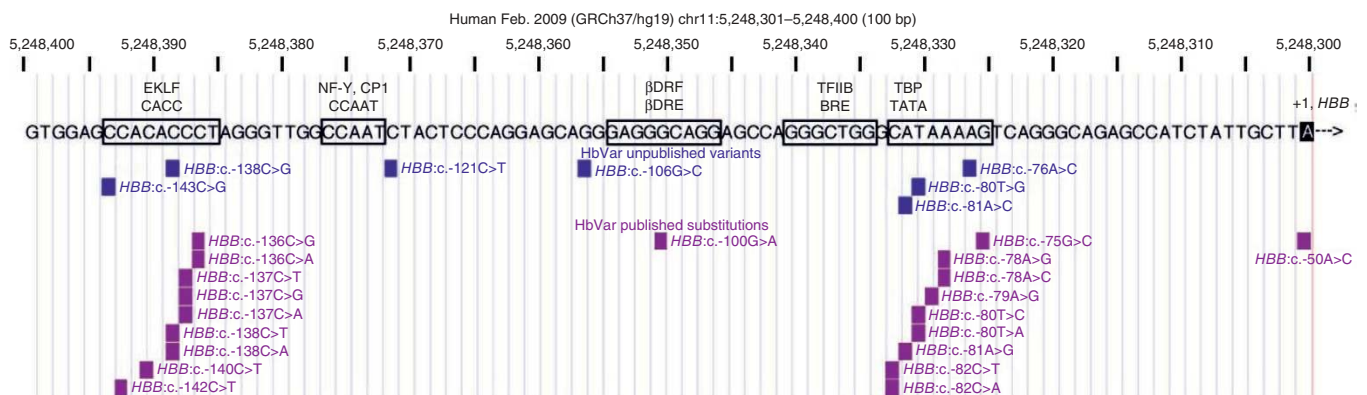
disease. Implementation of the microattribution approach has substantially added to the repository of variants, and use of this expanded database will continue to provide an important resource for generating and testing new hypotheses in the globin field. Below, we provide some recent examples illustrating the value of the microattribution approach. The value of the comprehensive globin variant database (pre- and post-microattribution) clearly emphasizes the importance of developing similar databases for other genes and disease systems for which microattribution will become the main route to publication.

The first example of the value of the microattribution approach is the finding that the distribution of promoter mutations differs among globin genes. Although a great deal has been learned about mammalian promoters from previous analyses of the globin genes, the discovery of additional variants continues to develop our knowledge of how these genes are normally activated and how they are altered in human genetic disease. Globin gene promoter mutations contributing to  $\beta$ -like thalassemias and HPFH comprise approximately 10% of the total variants and result in various phenotypes, from the asymptomatic non-deletional HPFH conditions to the mild forms of  $\beta$ - and  $\delta$ -thalassemia. The *HBB* promoter region harbors several genetic variants associated with  $\beta^+$  (expressing lower than normal levels of  $\beta$ -globin) and  $\beta^0$  (expressing no  $\beta$ -globin) thalassemia; these variants cluster in *cis*-regulatory elements known to bind transcription factors (Fig. 1). Many of these variants have been published, but an increasing number of unpublished variants have been contributed to HbVar by investigators around the world. The unpublished variants provide a more complete view of the contribution of genetic variants to phenotypes. In this particular case, they reveal phenotypic consequences of variants in more positions of well-known transcription factor binding sites (the CACC box and the TATA box) and show that additional substitutions in other binding sites contribute to phenotype (for example, positions c.-80, c.-81 and c.-138). The *HBB* c.-121C>T transition is adjacent to the CCAAT box. This motif was recognized 30 years ago as a component of some promoters, but the newly reported mutation here is the first indication that genetic variation close to this motif affects *HBB* gene expression in humans.

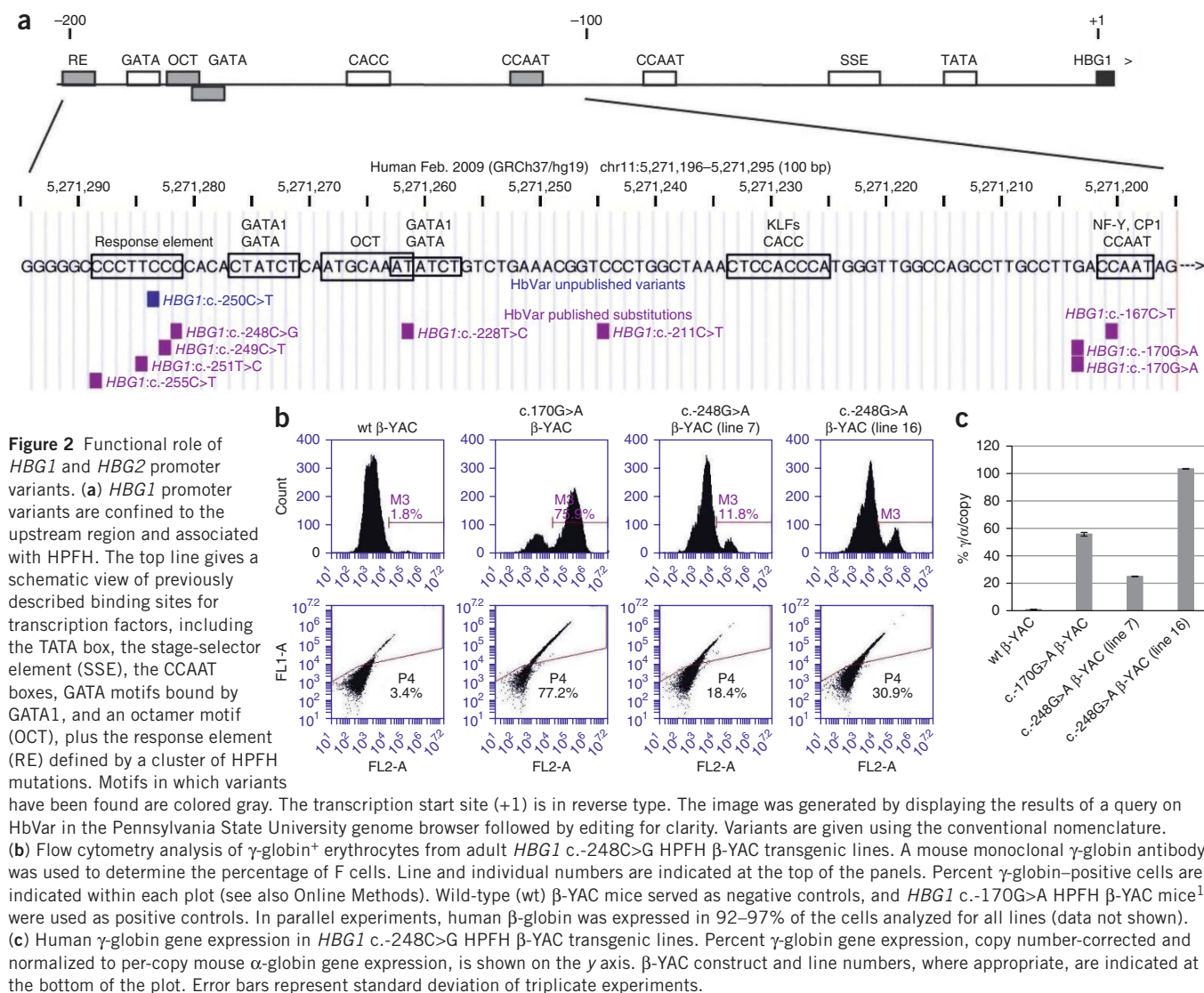
In contrast to the promoters for *HBB* and *HBD*, variants are not found in the first 100 bp of the *HBG1* and *HBG2* promoters, but instead, variants occur in the upstream region from approximately

-100 to -200 bp (Fig. 2a). The *HBG1* and *HBG2* gene promoters have several *cis*-regulatory elements in common with *HBB* and *HBD* promoters, such as a TATA box and a proximal CCAAT box, but no variants have been found in these elements. However, the CCAAT box is duplicated in the promoters of *HBG1* and *HBG2*, and the upstream CCAAT box (and the nucleotides very close to it) carries variants associated with HPFH. A newly discovered, unpublished variant, c.-250C>T, calls attention to a tight cluster of mutations all associated with HPFH. An HPFH-associated variant has now been reported at each nucleotide from position c.-251 to c.-248 (198 to 195 bp from the gene transcription start site), and a variant at c.-255 (202 bp from the transcription site) is associated with a similar phenotype (Fig. 2a). Given these phenotypes, this cluster of variants within the motif CCCTTCCC delineates a response element important for the silencing of the *HBG1* and presumably *HBG2* genes in adult erythroid cells (the same c.-250C>T mutation has been found in the promoter of *HBG2*; data not shown).

To test the hypothesis, derived from the documented variants, that this motif delineates a response element important for silencing of the *HBG1* and *HBG2* genes, we generated human  $\beta$ -globin locus ( $\beta$ -yeast artificial chromosome ( $\beta$ -YAC)) transgenic mice containing the *HBG1* c.-248C>G variation (the Brazilian non-deletional HPFH mutation), which directly alters the CCCTTCCC sequence at the 3' C. Adult mutant  $\beta$ -YAC mice showed an HPFH phenotype with an increased number of HbF-containing cells (Fig. 2b), and real-time quantitative RT-PCR analyses showed an 8- to 34-fold increase of *HBG1* gene expression relative to wild-type  $\beta$ -YAC mice (Fig. 2c). By comparison,  $\beta$ -YAC transgenic mice bearing the Greek type of non-deletional HPFH (*HBG1* c.-170G>A)<sup>13</sup> showed a 56-fold increase of *HBG1* gene expression relative to wild-type  $\beta$ -YAC mice. Future experiments will examine the mechanism of repression at this region. Recent studies have shown that the transcription factor BCL11A acts to repress *HBG1* and *HBG2* expression in adult erythroid cells, acting with the protein SOX6 (ref. 14). Although BCL11A showed no binding to the *HBG1* and *HBG2* proximal promoters, SOX6 showed strong binding that overlapped with GATA1 binding in these regions. In this way, the database has posed a new testable hypothesis. The CCCTTCCC element, which is adjacent to a GATA binding site, may bind a currently unknown protein that acts in concert with BCL11A to repress the production of  $\gamma$ -globins.



**Figure 1** Graphical display of the *HBB* promoter variants recorded in HbVar, partitioned into unpublished variants contributed by investigators (blue) and published variants (purple). The genomic position, sequence change and associated phenotype ( $\beta^+$  or  $\beta^0$  thalassemia) are given for each variant. Known protein-binding sites in the DNA sequence are boxed, with the name of the site and the binding protein above it. The transcription start site (+1) is in reverse type. The reverse complement of the genomic sequence is shown so that the gene is in the conventional left-to-right transcriptional orientation. The image was generated by displaying the results of a query on HbVar in the Pennsylvania State University genome browser followed by editing for clarity. Variants are given using the conventional nomenclature.



Overall, comparative analysis of the globin gene promoter mutations revealed a distinct distribution pattern for each gene. In *HBD*, promoter mutations are widely spread within the proximal promoter region and do not form mutational clusters around *cis*-regulatory elements (Supplementary Fig. 3). Notably, the mutations c.-81A>G and c.-80T>C have been found in the TATA boxes of *HBB* and *HBD*, suggesting that they could be the result of genetic recombination events<sup>15</sup>.

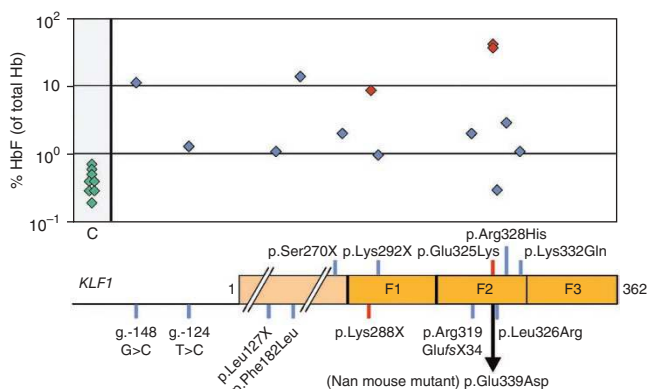
A second example of the value of the microattribution approach was the discovery of  $\alpha$ -thalassemia resulting from inherited or acquired mutations in *ATRX*. The comprehensive database originally identified and defined some of the key *trans*-acting factors in the globin gene system. The expanded database continues to refine our understanding of such *trans*-acting factors. Unlike the common forms of  $\alpha$ -thalassemia resulting from *cis*-acting genetic defects, two rare forms of  $\alpha$ -thalassemia are caused by *trans*-acting mutations in the X-linked *ATRX*. These mutations cause ATR-X syndrome, which is characterized by a severe form of syndromal mental retardation with characteristic dysmorphic faces, genital abnormalities and a mild but variable form of hemoglobin H disease<sup>3</sup>. In addition, acquired mutations in *ATRX* are seen in individuals who develop ATMDS syndrome,

a condition in which  $\alpha$ -thalassemia (AT) is associated with myelodysplastic syndrome (MDS)<sup>4</sup>. In both conditions, the levels of  $\alpha$ -globin mRNA are reduced, suggesting that *ATRX* is involved in the normal regulation of  $\alpha$ -globin gene expression. To date, 107 unique inherited and/or acquired disease-causing missense mutations have been found, which are located predominantly in two highly conserved domains of *ATRX* (Supplementary Fig. 4). These variants cluster within a globular domain that contains a plant homeodomain, which binds the N-terminal tails of histone H3, and the 7 helicase sub-domains, which identify *ATRX* as a member of the SNF2 family of chromatin-associated proteins. Structure and function studies based on natural mutations in the comprehensive database have elucidated precisely how *ATRX* is recruited to some of its targets through an interaction with the N-terminal tails of histone H3.

Notably, the degree of  $\alpha$ -thalassemia seen in individuals with ATMDS (having acquired *ATRX* gene mutations) is much greater than in individuals with the ATR-X syndrome (having inherited *ATRX* gene mutations), even when, by comparing mutations on the comprehensive database, we can see that the same *ATRX* mutation occurs in both conditions<sup>16</sup>. Again, analysis of the comprehensive variant database poses a new testable hypothesis. These findings suggest that

another component of the ATRX pathway may frequently be mutated in individuals with the common forms of MDS.

A third example of the value of microattribution is the discovery of variants in *KLF1* leading to elevated HbF levels. *KLF1* encodes a key erythroid transcriptional regulator that has many target genes with essential functions in erythroid cells including the globins, membrane proteins and heme synthesis enzymes<sup>17</sup>. The first report on *KLF1* mutations in humans linked them to the rare blood group In(Lu) phenotype<sup>18</sup>, in which the expression of the Lutheran blood group antigens is diminished. The reported individuals carried eight different loss-of-function mutations and one mutation abolishing a GATA1 binding site in the *KLF1* promoter. In all cases, the mutant *KLF1* allele occurred in the presence of a normal *KLF1* allele. A subsequent study on a large Maltese pedigree demonstrated that haploinsufficiency for *KLF1* causes HPPH<sup>7</sup>. A mutation in *KLF1*, resulting in p.Lys288X, was present exclusively in all individuals in this family with HPPH. This mutation ablates the complete zinc finger domain and therefore abrogates DNA binding of the mutant *KLF1* protein (Fig. 3 and Supplementary Table 2). The occurrence of HPPH in the individuals with In(Lu) has not been investigated. An analysis of archived blood samples from a number of these individuals with In(Lu) showed that their HbF levels were raised compared to those observed in control samples. Also, 30 out of 31 Sardinian individuals bearing four different *KLF1* mutations showed raised HbF levels compared to control samples. In addition, two individuals suffering from dyserythropoietic anemia carried a *KLF1* p.Glu325Lys alteration and had an HbF level of 40% (Fig. 3 and Supplementary Table 2)<sup>19,20</sup>. Mutations at this position alter the DNA binding specificity of *KLF1*. We note that the mouse neonatal anemia mutant (Nan) has an alteration in the orthologous amino acid of Klf1, p.Glu339Asp<sup>21,22</sup>. Adult heterozygous Nan animals show increased expression of embryonic globins, a condition akin to HPPH. Collectively, these data support the link between *KLF1* and HPPH and highlight the importance of the second DNA-binding zinc finger for normal *KLF1* function. This raises the possibility that some of the *KLF1* mutations which result in altered DNA binding specificity may have increased impact on HbF levels. This hypothesis can now be experimentally tested *in vitro* by DNA binding assays and *in vivo* in animal models.



**Figure 3** Correlation of the different *KLF1* gene variants deposited into HbVar (shown as blue and red squares, depicting unpublished and published information, respectively) and their corresponding HbF levels (median value in cases of three or more individuals) compared to wild-type individuals (shown as green squares). *KLF1* is not shown to scale. A simplified diagram depicting the *KLF1* promoter and protein is shown underneath. The positions of the zinc fingers are indicated (F1, F2 and F3). For the exact HbF levels corresponding to each *KLF1* gene variant, see Supplementary Table 2.

A final example of the value of microattribution is the discovery of hemoglobin variants. A large proportion of genetic variation in the human globin genes leads to hemoglobin variants. Most hemoglobin variants are rare, result from single amino acid substitutions of a globin chain and have a negligible or even no effect on hemoglobin function<sup>2</sup>.

The documented hemoglobin variants reside solely within exons and include: (i) structural variants with a pleiotropic effect (for example, HbS (*HBB* c.20A>T), HbE (*HBB* c.79G>A) and HbC (*HBB* c.19G>A)); (ii) variants (138 different variants) leading to unstable hemoglobin, where mutations affect the heme pocket of the globin chain; (iii) variants leading to methemoglobinemia, where the ferrous ion (Fe<sup>2+</sup>) of the heme group is oxidized to the ferric state (Fe<sup>3+</sup>) (most of these variants involve replacement by tyrosine of the histidine residues that anchor heme); and (iv) variants (92 different variants) with altered oxygen affinity, most of which result in increased oxygen affinity.

Although all of these correlations between structure and function have depended on data from the comprehensive database, new insights and questions continue to arise as new mutants are added to the repository, an initiative that sparked the implementation of the microattribution process for hemoglobinopathies. Notably, 14 hemoglobin variants result from the same mutation, but this mutation occurs on a different  $\alpha$ -globin gene paralogue<sup>23</sup>, that is, variations involving related genes that have evolved from recent gene duplication and as such are subject to frequent gene conversion events (Supplementary Fig. 5). HbF-Sardinia and HbF-Lesvos provide another such example, involving the same mutation (c.227T>C) but on the paralogous *HBB1* and *HBB2* genes, respectively<sup>24</sup>.

## DISCUSSION

The development of an integrated set of comprehensive LSDBs for a particular spectrum of human genetic diseases with microattribution, as described here for the hemoglobinopathies, provides an example of how such systems might be set up for a wide range of human genetic disorders in the future. Using the microattribution process set out here, datasets which took decades to accumulate for the globin genes could be assembled rapidly for other genes and disease systems. In the past, the description of natural variants has been accommodated by the conventional literature and has made an enormous contribution to the field of human genetics. In addition, it has shown how some of these mutations have reached polymorphic frequencies through natural selection, and detailed analysis of natural mutants has also been invaluable in establishing many of the general principles underlying mammalian gene regulation and human molecular genetics.

The strength of such observations will continue to increase as new mutations enter the databases, even though these might not merit a full publication on their own. Furthermore, new patterns of mutation may emerge; the accumulation of coding mutations in particular regions of a protein often identify a functionally important domain, as illustrated by *ATRX* and *KLF1* gene variants (Supplementary Fig. 4 and Fig. 3, respectively), and conversely, the identification of common neutral variants may rule out a major functional role for other regions. Similarly, DNA variants of key regulatory regions (promoters, enhancers, silencers, boundary elements and locus control regions) are often critical in identifying important *cis* elements and yet other neutral variants may help map regions of little functional importance (Fig. 1 and Supplementary Fig. 3). At the nucleotide level, such variants can even help map transcription factor binding sites<sup>25</sup>. The emergence of patterns of mutation may also point to the mechanisms of mutation, exemplified by gene conversion events identified at the

*HBA1* and *HBA2*, and *HBG1* and *HBG2* genes. Additionally, subtle phenotypic differences, for example, between  $\delta\beta$ -thalassemia and deletional HPFH<sup>2</sup>, can be attributed to the different junction points and the sequences that are removed or juxtaposed as a result of these deletions. Systematic documentation of these deletions in HbVar is currently under way and may allow for the identification of new regulatory elements that lie within the deleted or juxtaposed regions.

Perhaps the most important aspect of such comprehensive interacting databases is that they will pose and answer questions that would otherwise not be addressed, potentially leading to useful new insights. These databases will not only be of value in establishing the phenotypes of natural variants but may also be used in the development of personalized medicine. In the globin field, a great deal of effort is directed toward the development of drugs to increase the level of HbF and thereby ameliorate the clinical severity of  $\beta$ -thalassemia and SCD. Potential therapeutic agents identified to date include hydroxyurea and butyrate. The response to HbF-augmenting therapies is variable in patients with  $\beta$ -thalassemia and SCD, with approximately 25% of these patients being poor responders or non-responders<sup>26</sup>. Therefore, the ability to predict a patient's response to hydroxyurea and/or other HbF-augmenting drugs would help in optimizing therapy. Polymorphisms in genes regulating HbF expression, hydroxyurea metabolism and erythroid progenitor proliferation might modulate a patient's response to HbF-inducing pharmacological agents<sup>27</sup>. Data to support the use of pharmacogenetic testing of hydroxyurea treatment for hemoglobinopathies are currently very limited. Several SNPs in *HAO2*, *ARG2*, *FLT1* and *NOS1* have been associated with variable HbF response to hydroxyurea treatment<sup>27</sup>, and genome-wide transcription profiling efforts are expected to shed light on new pathways involved in this process<sup>28</sup>.

Since its establishment in 2000, we have witnessed a substantial annual growth in HbVar content, and a fraction of data submitters were subsequently encouraged to submit a full or short report to the scientific journal *Hemoglobin*<sup>29</sup>. The large repository of previously reported data, together with more recent data acquired by microattribution, shows how the comprehensive documentation of human variation will provide key insights into normal biological processes and how these are perturbed in human genetic disease. We anticipate that microattribution will further encourage new data submitters to contribute their observations to HbVar to receive not only credit in the form of microcitations but also coauthorship in a future microattribution update. The microattribution process established here provides a template for similar ventures for other human genes, their associated systems and the variants that cause their associated genetic diseases. The value of the databases may be considerably further enhanced by linking to collections of blood and DNA samples and also cataloged online, as in the case of many other rare diseases in EuroBioBank.

In essence, this project is a well-coordinated multicenter effort to systematically document genetic variation in globin and associated genes relevant to hemoglobinopathies and thalassemias and is the first example of implementing microattribution to provide incentives for submitting data describing genetic variation. As such, it should serve as a model for the comprehensive documentation and analysis of genetic variations in other common or genetically complex disorders, the conduct of a thorough synopsis of other fields, or both.

**URLs.** HbVar Database of Hemoglobin Variants and Thalassemia Mutations, <http://globin.bx.psu.edu/hbvar/>; Golden Helix Server, <http://www.goldenhelix.org/>; Leiden Open-Access Variation Database, <http://www.lov.d.nl/>; Frequencies of Inherited Disorders database, <http://www.findbase.org/>; dbSNP database, <http://www.ncbi.nlm.nih.gov/projects/SNP/>; Human Genome Variation Society, <http://www.hgvs.org/>; ResearcherID System of Thomson ISI, <http://www.researcherid.com/>; Open ID system, <http://openid.net/>; Genotype-to-Phenotype database project's Researcher Identification Primer (RIP), <http://www.gen2phen.org/>.

nih.gov/projects/SNP/; Human Genome Variation Society, <http://www.hgvs.org/>; ResearcherID System of Thomson ISI, <http://www.researcherid.com/>; Open ID system, <http://openid.net/>; Genotype-to-Phenotype database project's Researcher Identification Primer (RIP), <http://www.gen2phen.org/>.

## METHODS

Methods and any associated references are available in the online version of the paper at <http://www.nature.com/naturegenetics/>.

*Note: Supplementary information is available on the Nature Genetics website.*

## ACKNOWLEDGMENTS

This work was supported by the Netherlands Genomics Initiative (NGI), Erasmus MC (MRace; 296088), the Landsteiner Foundation for Blood Transfusion Research (LSBR; 1040), US National Institutes of Health (NIH) (R01-HL073455) and the Netherlands Scientific Organization (NWO DN 82-301 and 912-07-019) to S.P., the NIH grants R01 DK065806, RC HG005573 and U01 HG004695 to R.C.H., the National Institutes for Health Research Biomedical Research Centre (Oxford) to D.R.H. and European Commission grants (FP6-026539 (ITHANET), FP7-200754 (GEN2PHEN)) to G.P.P.

## AUTHOR CONTRIBUTIONS

B.G., J.B., R.C.H. and G.P.P. conceived and designed the study. B.G., J.B. and D.M. implemented the process, built and populated the databases. K.R.P., F.C.C. and H.F. performed experiments. B.K.S., D.J.A., A.N.B., B.C., P.F., A.E.F., A.F., R.G., M.V.E.G., M.G., R.J.G., P.C.G., C.L.H., J.D.H., M.J., P.J., E.K., P.K., S.M., K.M., J.O., A.P., M.N.P., P.P., S. Pavlovic, L.P., M.R., S.S., I.S., M.S., S.L.T., J.T.-S., R.T., T.W., J.S.W., C.W., B.Z. and G.P.P. contributed data. B.G., J.B., D.R.H. and S. Philipsen analyzed results. R.C.H. and G.P.P. supervised data analysis. D.M., W.M., C.R., D.H.K.C. and H.W. provided expertise and infrastructure. B.G., J.B., D.R.H., K.R.P., S. Philipsen, R.C.H. and G.P.P. wrote the paper.

## COMPETING FINANCIAL INTERESTS

The authors declare no competing financial interests.

Published online at <http://www.nature.com/naturegenetics/>.

Reprints and permissions information is available online at <http://npg.nature.com/reprintsandpermissions/>.

This paper is distributed under the terms of the Creative Commons Attribution-Noncommercial-Share Alike license, and is freely available to all readers at <http://www.nature.com/naturegenetics/>.

1. Anonymous. Human variome microattribution reviews. *Nat. Genet.* **40**, 1 (2008).
2. Patrinos, G.P. & Antonarakis, S.E. Human hemoglobin. in *Human Genetics: Problems and Approaches* (eds Speicher, M., Antonarakis, S.E. & Motulsky, A.), 366–401 (Springer-Verlag, Heidelberg, Germany, 2010).
3. Gibbons, R.J., Picketts, D.J., Villard, L. & Higgs, D.R. Mutations in a putative global transcriptional regulator cause X-linked mental retardation with alpha-thalassemia (ATR-X syndrome). *Cell* **80**, 837–845 (1995).
4. Gibbons, R.J. *et al.* Identification of acquired somatic mutations in the gene encoding chromatin-remodeling factor ATRX in the  $\alpha$ -thalassemia myelodysplasia syndrome (ATMDS). *Nat. Genet.* **34**, 446–449 (2003).
5. Viprakasit, V. *et al.* Mutations in the general transcription factor TFIIF result in  $\beta$ -thalassaemia in individuals with trichothiodystrophy. *Hum. Mol. Genet.* **10**, 2797–2802 (2001).
6. Yu, C. *et al.* X-linked thrombocytopenia with thalassemia from a mutation in the amino finger of GATA-1 affecting DNA binding rather than FOG-1 interaction. *Blood* **100**, 2040–2045 (2002).
7. Borg, J. *et al.* Haploinsufficiency for the erythroid transcription factor KLF1 causes hereditary persistence of fetal hemoglobin. *Nat. Genet.* **42**, 801–805 (2010).
8. Thein, S.L. *et al.* Intergenic variants of *HBS1L-MYB* are responsible for a major quantitative trait locus on chromosome 6q23 influencing fetal hemoglobin levels in adults. *Proc. Natl. Acad. Sci. USA* **104**, 11346–11351 (2007).
9. Menzel, S. *et al.* A QTL influencing F cell production maps to a gene encoding a zinc-finger protein on chromosome 2p15. *Nat. Genet.* **39**, 1197–1199 (2007).
10. Sankaran, V.G. *et al.* Human fetal hemoglobin expression is regulated by the developmental stage-specific repressor BCL11A. *Science* **322**, 1839–1842 (2007).
11. Hardison, R.C. *et al.* HbVar: a relational database of human hemoglobin variants and thalassemia mutations at the globin gene server. *Hum. Mutat.* **19**, 225–233 (2002).
12. Fokkema, I.F., den Dunnen, J.T. & Taschner, P.E. LOVD: easy creation of a locus-specific sequence variation database using an "LSDB-in-a-box" approach. *Hum. Mutat.* **26**, 63–68 (2005).

13. Peterson, K.R. *et al.* Use of yeast artificial chromosomes (YACs) in studies of mammalian development: production of  $\beta$ -globin locus YAC mice carrying human globin developmental mutants. *Proc. Natl. Acad. Sci. USA* **92**, 5655–5659 (1995).
14. Xu, J. *et al.* Transcriptional silencing of  $\gamma$ -globin by BCL11A involves long-range interactions and cooperation with SOX6. *Genes Dev.* **24**, 783–798 (2010).
15. Borg, J., Georgitsi, M., Aleporou-Marinou, V., Kollia, P. & Patrinos, G.P. Genetic recombination as a major cause of mutagenesis in the human globin gene clusters. *Clin. Biochem.* **42**, 1839–1850 (2009).
16. Steensma, D.P., Gibbons, R.J. & Higgs, D.R. Acquired  $\alpha$ -thalassemia in association with myelodysplastic syndrome and other hematologic malignancies. *Blood* **105**, 443–452 (2005).
17. Drissen, R. *et al.* The erythroid phenotype of EKLF-null mice: defects in hemoglobin metabolism and membrane stability. *Mol. Cell. Biol.* **25**, 5205–5214 (2005).
18. Singleton, B.K., Burton, N.M., Green, C., Brady, R.L. & Anstee, D.J. Mutations in EKLF/KLF1 form the molecular basis of the rare blood group In(Lu) phenotype. *Blood* **112**, 2081–2088 (2008).
19. Arnaud, L. *et al.* A dominant mutation in the gene encoding the erythroid transcription factor KLF1 causes a congenital dyserythropoietic anemia. *Am. J. Hum. Genet.* **87**, 721–727 (2010).
20. Singleton, B.K. *et al.* A novel EKLF mutation in a patient with dyserythropoietic anemia: the first association of EKLF with disease in man. *Blood* **114**, 72 (2009).
21. Siatecka, M. *et al.* Severe anemia in the Nan mutant mouse caused by sequence-selective disruption of erythroid Kruppel-like factor. *Proc. Natl. Acad. Sci. USA* **107**, 15151–15156 (2010).
22. Heruth, D.P. *et al.* Mutation in erythroid specific transcription factor KLF1 causes hereditary spherocytosis in the Nan hemolytic anemia mouse model. *Genomics* **96**, 303–307 (2010).
23. Moradkhani, K. *et al.* Mutations in the paralogous human  $\alpha$ -globin genes yielding identical hemoglobin variants. *Ann. Hematol.* **88**, 535–543 (2009).
24. Papadakis, M.N., Patrinos, G.P., Drakoulakou, O. & Loutradi-Anagnostou, A. HbF-Lesvos: an HbF variant due to a novel G gamma mutation (:G gamma 75 ATA→ACA) detected in a Greek family. *Hum. Genet.* **97**, 260–262 (1996).
25. Patrinos, G.P. *et al.* Improvements in the HbVar database of human hemoglobin variants and thalassemia mutations for population and sequence variation studies. *Nucleic Acids Res.* **32**, D537–D541 (2004).
26. Steinberg, M.H. *et al.* Fetal hemoglobin in sickle cell anemia: determinants of response to hydroxyurea. Multicenter study of hydroxyurea. *Blood* **89**, 1078–1088 (1997).
27. Ma, Q. *et al.* Fetal hemoglobin in sickle cell anemia: genetic determinants of response to hydroxyurea. *Pharmacogenomics J.* **7**, 386–394 (2007).
28. Patrinos, G.P. & Grosveld, F.G. Pharmacogenomics and therapeutics of hemoglobinopathies. *Hemoglobin* **32**, 229–236 (2008).
29. Patrinos, G.P. & Wajcman, H. Recording human globin gene variation. *Hemoglobin* **28**, v–vii (2004).

<sup>1</sup>Pennsylvania State University, Center for Comparative Genomics and Bioinformatics, University Park, Philadelphia, Pennsylvania, USA. <sup>2</sup>Department of Applied Biomedical Sciences, University of Malta, Msida, Malta. <sup>3</sup>Laboratory of Molecular Genetics, Department of Physiology and Biochemistry, University of Malta, Msida, Malta. <sup>4</sup>Thalassemia Clinic, Section of Pathology, Mater Dei Hospital, Msida, Malta. <sup>5</sup>Medical Research Council (MRC) Molecular Haematology Unit, Weatherall Institute of Molecular Medicine, Oxford, UK. <sup>6</sup>University of Kansas Medical Center, Department of Biochemistry and Molecular Biology, Kansas City, Kansas, USA. <sup>7</sup>Erasmus University Medical Center, Faculty of Medicine and Health Sciences, Department of Cell Biology, Rotterdam, The Netherlands. <sup>8</sup>National Center for Biotechnology Information, National Library of Medicine, National Institutes of Health, Bethesda, Maryland, USA. <sup>9</sup>Bristol Institute for Transfusion Sciences (BITS), National Health Service (NHS) Blood and Transplant, Bristol, UK. <sup>10</sup>Bogazici University, Department of Molecular Biology and Genetics, Istanbul, Turkey. <sup>11</sup>King's College London, London, UK. <sup>12</sup>Unidade de Investigação e Desenvolvimento, Departamento de Genética, Instituto Nacional de Saúde Dr. Ricardo Jorge, Lisboa, Portugal. <sup>13</sup>Department of Biochemistry, Edouard Herriot University Hospital, Lyon Cedex, France. <sup>14</sup>Dipartimento di Scienze Biomediche e Biotecnologie, University of Cagliari, Cagliari, Sardinia, Italy. <sup>15</sup>Quest Diagnostics Nichols Institute, Chantilly, Virginia, USA. <sup>16</sup>Department of Pharmacy, School of Health Sciences, University of Patras, Patras, Greece. <sup>17</sup>Hemoglobinopathies Laboratory, Human and Clinical Genetics Department, Leiden University Medical Center, Leiden, The Netherlands. <sup>18</sup>Mayo Clinic, Division of Hematopathology, Rochester, Minnesota, USA. <sup>19</sup>North Middlesex University Hospital, London, UK. <sup>20</sup>National and Kapodistrian University of Athens, School of Medicine, Medical Genetics, St. Sophia's Children's Hospital, Athens, Greece. <sup>21</sup>Department of Biology, National and Kapodistrian University of Athens, School of Physical Sciences, Athens, Greece. <sup>22</sup>Hospital Henri-Mondor and Albert-Chenevier Group, Department of Biochemistry and Genetics, Créteil, France. <sup>23</sup>National Haemoglobinopathy Reference Laboratory, Oxford Haemophilia Centre, Churchill Hospital, Oxford, UK. <sup>24</sup>University of Patras, Faculty of Medicine, Laboratory of General Biology, Patras, Greece. <sup>25</sup>Unit of Prenatal Diagnosis, Center for Thalassemia, Laikon General Hospital, Athens, Greece. <sup>26</sup>Institute of Molecular Genetics and Genetic Engineering, University of Belgrade, Belgrade, Serbia. <sup>27</sup>Istituto di Neurogenetica e Neurofarmacologia, National Research Council, Cagliari, Cagliari, Sardinia, Italy. <sup>28</sup>Stanford University School of Medicine, Department of Pathology and Pediatrics, Stanford, California, USA. <sup>29</sup>Division of Neurology, Kanagawa Children's Medical Center, Yokohama, Kanagawa, Japan. <sup>30</sup>Department of Pathology and Molecular Medicine, McMaster University, Hamilton, Ontario, Canada. <sup>31</sup>Molecular Diagnostic Genetics, Hamilton Regional Laboratory Program, Hamilton, Ontario, Canada. <sup>32</sup>Medizinisches Versorgungszentrum (MVZ), Laboratory Prof. Seelig, Karlsruhe, Germany. <sup>33</sup>Department of Medicine, Boston University School of Medicine, Boston, Massachusetts, USA. <sup>34</sup>Department of Pathology, Boston University School of Medicine, Boston, Massachusetts, USA. <sup>35</sup>INSERM, U955, Créteil, France. <sup>36</sup>Department of Biochemistry and Molecular Biology, Pennsylvania State University, University Park, Philadelphia, Pennsylvania, USA. <sup>37</sup>These authors contributed equally to this work. Correspondence should be addressed to G.P.P. (gpatrinos@upatras.gr).

## ONLINE METHODS

**Quantitation of hemoglobin fractions.** Twenty microliters of total blood was analyzed using cation-exchange high performance liquid chromatography (VARIANT, Bio-Rad Laboratories).

**Construction of the *HBG1* c.-248C>G HPHF  $\beta$ -YAC.** A 213-kb yeast artificial chromosome (YAC) carrying the human  $\beta$ -globin locus with the *HBG1* c.-248C>G point mutation ( $^A\gamma$ -195C>G), leading to the Brazilian type of non-deletional HPHF, which directly alters the CCCTTCCC sequence at the 3' C, was synthesized as follows, using previously described methods<sup>30</sup>. Briefly, a marked *HBG1* gene ( $^A\gamma^m$ ) contained as a 5.4-kb *SspI* fragment (GenBank file U01317, coordinates 38,683-44,077) in the yeast-integrating plasmid (YIP) pRS406 was mutagenized using the QuikChange Site-Specific Mutagenesis Kit (Stratagene). The presence of the *HBG1* c.-248C>G point mutation was confirmed by DNA sequencing, and the mutation was introduced into the  $\beta$ -YAC by 'pop-in', 'pop-out' homologous recombination in yeast. The mark in the  $^A\gamma^m$ -globin gene is a 6-bp deletion at +21 to +26 relative to the  $^A\gamma$ -globin translation start site, allowing preliminary discrimination of the modified  $\beta$ -YAC from the wild-type  $\beta$ -YAC by restriction enzyme digestion following homologous recombination. The presence of the mutation in clones passing this test was confirmed by DNA sequence analysis of a PCR-amplified fragment encompassing the mutated region. Transformation of yeast, screening of positive clones, purification of the  $\beta$ -YAC and mouse transgenesis were performed as described previously<sup>31</sup>.

**Copy number determination.** The relative  $\beta$ -YAC transgene copy number was calculated using the *HBG1* and *HBG2* genes and a standard curve generated from genomic DNA samples from our wild-type  $\beta$ -YAC transgenic mice. Samples of transgenic mouse genomic DNA were serially diluted from 100–0.01 ng and subjected to SYBR PCR with *HBG1* or *HBG2* primers. The copy number for each reaction was estimated by comparing the threshold cycle of each sample to the threshold cycle of the standards and normalizing to the wild-type  $\beta$ -YAC transgenic mouse samples.

**Real-time quantitative RT-PCR.** Total RNA, isolated from adult peripheral blood, was reverse-transcribed and the resultant complementary DNA was subjected to real-time quantitative RT-PCR analysis with SYBR green using a CFX96 system (Bio-Rad). Human  $\gamma$ -globin expression was normalized to mouse  $\alpha$ -globin expression and corrected for transgene and endogenous gene copy number. PCR primer sequences were as previously described<sup>32</sup>. Results are averages of triplicates, with the standard error indicated.

**F-cell detection by flow cytometry.** We used a protocol adapted from references 32 and 33. Essentially, mouse blood was collected from the tail vein in heparinized capillary tubes. Ten microliters of whole blood was washed in 1 ml PBS, centrifuged at 200g at 4 °C for five minutes, and the pellet was resuspended and fixed in 1 ml of 4% fresh paraformaldehyde and PBS at pH 7.5 (Sigma-Aldrich) for 40 min at 37 °C. The cells were centrifuged, and the pellets were resuspended in 1 ml of ice cold acetone and methanol (4:1) and incubated on ice for one minute. Following centrifugation, cells were washed twice in 1 ml ice-cold PBS and 0.1% BSA and resuspended in 800  $\mu$ l of PBS, 0.1% BSA and 0.1% Triton X-100 (PBT). One microgram of  $\gamma$ -globin antibody (catalog number sc-21756 unconjugated, Santa Cruz Biotechnology) was added to 100  $\mu$ l of the cell suspension and incubated for 20 min in the dark at room temperature (37 °C). One milliliter of ice-cold PBS and 0.1% BSA was added, the sample was centrifuged and the pellet was resuspended in 100  $\mu$ l ice-cold PBT. One hundred microliters of Alexa 488 (catalog number 11001, Invitrogen Molecular Probes) secondary antibody, diluted 1:200 in ice-cold PBT, was added to the cell suspension and the sample was incubated at room temperature for 20 min in the dark. Cells were washed with 1 ml of ice-cold PBS and 0.1% BSA and the pellets were resuspended in 200  $\mu$ l of PBS. Samples were analyzed using an Accuri C6 Flow Cytometer (Accuri Cytometers, Inc.) with a 530/30 nm (FITC/GFP) emission filter. Data from 30,000 cells were acquired for analysis using CFlow Software (Accuri Cytometers, Inc.); cells were gated to exclude dead cells. For, FL1-A, a 530/30 nm (FITC, GFP) filter was used to identify the Alexa 488-positive F cell population; For FL2-A a 585/40 nm (PE, PI) filter was used as a compensation to identify the Alexa 488-negative cell population. For M3, the mean fluorescent intensity, an increase in F cells is reflected by a peak shift and increase in the peak of fluorescence intensity. P4, distinct positive F cells.

30. Harju, S., Navas, P.A., Stamatoyannopoulos, G. & Peterson, K.R. Genome architecture of the human  $\beta$ -globin locus affects developmental regulation of gene expression. *Mol. Cell. Biol.* **25**, 8765–8778 (2005).
31. Harju-Baker, S., Costa, F.C., Fedosyuk, H., Neades, R. & Peterson, K.R. Silencing of Agamma-globin gene expression during adult definitive erythropoiesis mediated by GATA-1-FOG-1-Mi2 Complex binding at the -566 GATA site. *Mol. Cell. Biol.* **28**, 3101–3113 (2008).
32. Böhmer, R.M. Flow cytometry of erythroid cells in culture: bivariate profiles of fetal and adult hemoglobins. *Methods Cell Biol.* **64**, 139–152 (2001).
33. Amoyal, I. & Fibach, E. Flow cytometric analysis of fetal hemoglobin in erythroid precursors of  $\beta$ -thalassemia. *Clin. Lab. Haematol.* **26**, 187–193 (2004).