



# Identifying Recent Adaptations in Large-Scale Genomic Data

Sharon R. Grossman,<sup>1,2,4,14,\*</sup> Kristian G. Andersen,<sup>1,6,14</sup> Ilya Shlyakhter,<sup>1,6,14</sup> Shervin Tabrizi,<sup>1,6,14</sup> Sarah Winnicki,<sup>1,6</sup> Angela Yen,<sup>1,3</sup> Daniel J. Park,<sup>1,6</sup> Dustin Griesemer,<sup>4,6</sup> Elinor K. Karlsson,<sup>1,6</sup> Sunny H. Wong,<sup>8</sup> Moran Cabili,<sup>1,5</sup> Richard A. Adegbola,<sup>9</sup> Rameshwar N.K. Bamezai,<sup>10</sup> Adrian V.S. Hill,<sup>8</sup> Fredrik O. Vannberg,<sup>11</sup> John L. Rinn,<sup>1,7,12</sup> 1000 Genomes Project, Eric S. Lander,<sup>1,2,5</sup> Stephen F. Schaffner,<sup>1</sup> and Pardis C. Sabeti<sup>1,6,13,\*</sup>

<sup>1</sup>Broad Institute of MIT and Harvard, Cambridge, MA 02142, USA

<sup>2</sup>Department of Biology

<sup>3</sup>Computer Science and Artificial Intelligence Laboratory  
MIT, Cambridge, MA 02139, USA

<sup>4</sup>Division of Health Science and Technology

<sup>5</sup>Department of Systems Biology  
Harvard Medical School, Boston, MA 02115, USA

<sup>6</sup>Center for Systems Biology, Department of Organismic and Evolutionary Biology

<sup>7</sup>Stem Cell and Regenerative Biology  
Harvard University, Cambridge, MA 02138, USA

<sup>8</sup>Wellcome Trust Centre for Human Genetics, University of Oxford, Oxford OX3 7BN, UK

<sup>9</sup>MRC Laboratories, Fajara, the Gambia

<sup>10</sup>National Centre of Applied Human Genetics, School of Life Sciences, Jawaharlal Nehru University, New Delhi 110067, India

<sup>11</sup>School of Biology, Center for Integrative Genomics, Georgia Institute of Technology, Atlanta, GA 30332, USA

<sup>12</sup>Beth Israel Deaconess Hospital, Boston, MA 02115, USA

<sup>13</sup>Department of Immunology and Infectious Diseases, Harvard School of Public Health, Cambridge, MA 02115, USA

<sup>14</sup>These authors contributed equally to this work

\*Correspondence: grossman@broadinstitute.org (S.R.G.), psabeti@oeb.harvard.edu (P.C.S.)

<http://dx.doi.org/10.1016/j.cell.2013.01.035>

## SUMMARY

Although several hundred regions of the human genome harbor signals of positive natural selection, few of the relevant adaptive traits and variants have been elucidated. Using full-genome sequence variation from the 1000 Genomes (1000G) Project and the composite of multiple signals (CMS) test, we investigated 412 candidate signals and leveraged functional annotation, protein structure modeling, epigenetics, and association studies to identify and extensively annotate candidate causal variants. The resulting catalog provides a tractable list for experimental follow-up; it includes 35 high-scoring nonsynonymous variants, 59 variants associated with expression levels of a nearby coding gene or lincRNA, and numerous variants associated with susceptibility to infectious disease and other phenotypes. We experimentally characterized one candidate nonsynonymous variant in Toll-like receptor 5 (TLR5) and show that it leads to altered NF-κB signaling in response to bacterial flagellin.

## INTRODUCTION

Within their recent evolutionary history, humans traveled out of Africa into a wide range of new environments, endured repeated

climate change, and experienced dramatic alterations in diet and disease risk, exposing them to new and powerful forces. Although many recent selective pressures have been hypothesized, only a handful of adaptive traits have ever been characterized, such as malaria resistance in *HBB* (Currat et al., 2002; Ohashi et al., 2004), lactose tolerance in *LCT* (Bersaglieri et al., 2004), skin pigmentation in *SLC24A5* (Lamason et al., 2005), and high-altitude tolerance in *EPAS1* (Yi et al., 2010). Each of these began with knowledge of a phenotypic trait that was hypothesized to be adaptive and subsequently genetic evidence for selection was discovered.

The advent of genomics holds great promise for the study of human evolution, making it possible to move from hypothesis-driven candidate gene studies to hypothesis-generating genome-wide scans. Over the last decade, genome-wide scans for selection have been frequently reported, finding several hundred loci that show patterns of variation characteristic of new beneficial mutations that have spread quickly through the population (Akey, 2009; Akey et al., 2002; Bustamante et al., 2005; Frazer et al., 2007; Pickrell et al., 2009; Sabeti et al., 2007; Voight et al., 2006; Williamson et al., 2007).

Moving from candidate genomic regions, however, to the underlying adaptive mutation has been difficult. The dearth of adaptive variants elucidated from genomic scans has led some to question whether many such variants remain to be identified. For example, a study of diversity and differentiation in 1000 Genomes (1000G) Project data estimated that ~0.5% of nonsynonymous substitutions in the past 250,000 years have been subject to positive selection and concluded that adaptive

substitutions were rare in recent human history (Hernandez et al., 2011). However, this bound still allows for 340 nonsynonymous adaptive mutations in the 1000G data as well as countless regulatory mutations. The main reason so few examples have been identified is instead the difficulty of the problem. The regions detected are large, spanning hundreds of kilobases to megabases and containing thousands of potential variants driving the signal, whereas the relevant phenotypic trait is unknown. Furthermore, full-sequence data, necessary to compile a complete list of candidates, has not been available.

With this challenge in mind, we previously developed the method, the composite of multiple signals (CMS), designed to pinpoint a small number of candidate selected variants within a large genomic region (Grossman et al., 2010). CMS combines several independent population genetic statistics that distinguish the causal variant from neighboring neutral variants, thus reducing the number of candidates by between 20- and 100-fold while maintaining high sensitivity for the causal variant. In that study, we applied it as a proof-of-concept to simulations and to published candidate genomic regions from the International Haplotype Map (HapMap) Project, but the work was fundamentally limited in its ability to find underlying causal variants by incomplete genotype data. One obtains far greater power to narrow regions with full-sequence data. Furthermore, the sequence data ensure that all potential candidate variants are examined before pursuing functional validation.

Here, we present a catalog of candidate selected mutations, rather than genomic regions, using full-genome sequencing data from the 1000G Project. The result is a tractable set of variants for follow-up functional characterization. Initial functional annotation reveals coding mutations as well as many variants in regulatory regions and noncoding RNAs and associations with a variety of phenotypes. We use this database to identify a nonsynonymous mutation in *TLR5* with strong evidence for selection and show it leads to altered NF- $\kappa$ B signaling in response to stimulation with bacterial flagellin. This example, together with another paper in this issue characterizing a selected variant in the Ectodysplasin receptor (*EDAR*), present a framework for moving from a genome-wide scan, to a tractable set of candidate variants, to insights from functional annotation, to characterization of a population-specific functional variant, and elucidation of distinct mechanisms of human evolutionary adaptation.

## RESULTS

### 412 Fine-Mapped Signals of Selection

The deep characterization of human sequence variation produced by the 1000G Project permits examination of the vast majority of nucleotides in the genome for evidence of recent evolution. The pilot phase of 1000G included complete, low-coverage (2–6×), whole-genome sequencing of 179 individuals from four populations: Yoruba individuals from Nigeria (YRI), European-ancestry individuals from Utah (CEU), Han Chinese individuals in Beijing (CHB) and Japanese individuals in Tokyo (JPT) (1000 Genomes Project Consortium, 2010). The resulting data set covers approximately 85% of the reference sequence and 93% of the coding sequence of the genome, with the vast

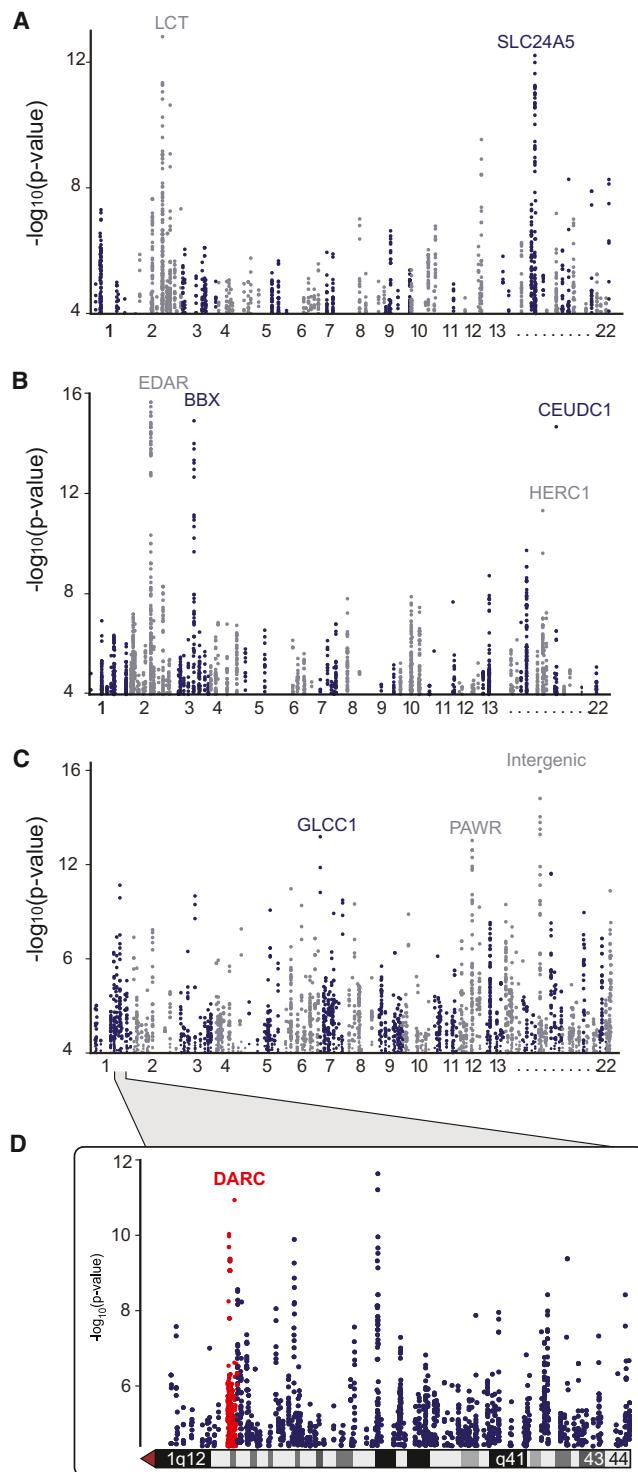
majority (97%) of inaccessible sites being high-copy repeats or segmental duplications. Comparison with overlapping samples in independent studies indicates that the methodology used in the 1000G study has 90% power to detect single nucleotide polymorphisms (SNPs) present at 5% frequency in a population and 99% power for SNPs at 10% frequency or greater. The project also identified numerous insertion/deletion polymorphisms (indels). Because standard current methods primarily identify selective sweeps with causal mutations at high frequency (>20%) in at least one population, the causal variants for detectable loci under selection are almost certainly present in the 1000G sequence data.

We recently developed the CMS method to fine-map signals of natural selection within previously identified candidate regions (Grossman et al., 2010); however, we subsequently hypothesized that the method could be modified to also detect new candidate regions. We thus developed a genome-wide CMS method (CMS<sub>GW</sub>, Experimental Procedures) and found that it was complementary to long-haplotype methods used to identify published candidate regions in the HapMap Project. Long-haplotype signals (generated as the selected mutation rapidly rises in frequency, bringing neighboring variants along with it) are very powerful for detecting recent sweeps (<30,000 years ago), but their power falls off significantly for older events due to haplotype breakdown. By incorporating signals that persist longer, such as population differentiation and high-frequency derived alleles, CMS<sub>GW</sub> is able to capture older events, for example selection at the *DARC* locus in West Africa, ca. 50,000 years ago.

We used CMS<sub>GW</sub> to identify 86 regions likely to be under selection at an FDR of 19% (Figures 1 and S1 and Table S1 available online) and combined these with regions previously identified using long-haplotype methods (Frazer et al., 2007). Both kinds of tests have better power to detect incomplete sweeps (where the causal allele is not yet fixed in the population) compared with frequency-spectrum-based statistics such as Tajima's D and Fay and Wu's H. They are thus well suited to our focus on recent (<50,000 years old) human selection events for which the beneficial alleles are unlikely to have reached fixation.

With this combined set of candidate genomic regions in hand, we used our standard CMS implementation to fine-map signals and identify the candidate causal mutations within full-sequence data (Grossman et al., 2010). The output was a set of 20–100 candidate variants per region (median = 47) at a threshold that captured 90% of causal variants in simulations. The candidates lay within genomic regions with a median size of 27 kb (Table S2). As current functional annotation is incomplete, we included all candidate causal variants in our database, regardless of any prior knowledge of functionality.

A possible source of artifactual signals of selection is copy number variants (CNVs), which have been suggested to play a role in creating unusually long haplotypes (Gusev et al., 2012). However, we identified only 60 instances of overlap between the localized regions and CNVs, not significantly higher than the number of overlaps expected at random (50, p = 0.23). Although our analysis below focuses on SNPs, we cataloged these CNVs as they could themselves be targets of selection (Extended Experimental Procedures).



**Figure 1. Genome-wide CMS**

(A–D) CMS<sub>GW</sub> scores calculated from full-genome sequence data from the 1000G Consortium in samples from (A) Northern Europe, (B) East Asia, and (C) West Africa. (D) Close-up of CMS<sub>GW</sub> scores from West Africa in region on chromosome 1 containing DARC. The y axis represents the significance level ( $p$  value on  $-\log_{10}$  scale) for each of the variants tested across the genome, showing only variants with significance levels exceeding  $10^{-4}$  (corresponding to 4 on the y axis). See also Figure S1 and Table S1.

Among our fine-mapped candidates, 147 regions contain a single protein-coding gene, 88 regions contain multiple genes and 177 regions do not have any genes coding for known proteins. The regions are enriched for genes ( $p = 0.08$ ) and coding variants ( $p < 0.01$ ) and contain a number of genes involved in biological pathways thought to be recently targeted by selection, such as skin pigmentation and the immune system (Tables S3 and S4). They also contain 48 long intergenic noncoding RNAs (lincRNAs) (Cabili et al., 2011), thirteen of which lie in regions with no protein-coding genes, suggesting another class of functional elements that may be a target of recent positive selection (Table S5).

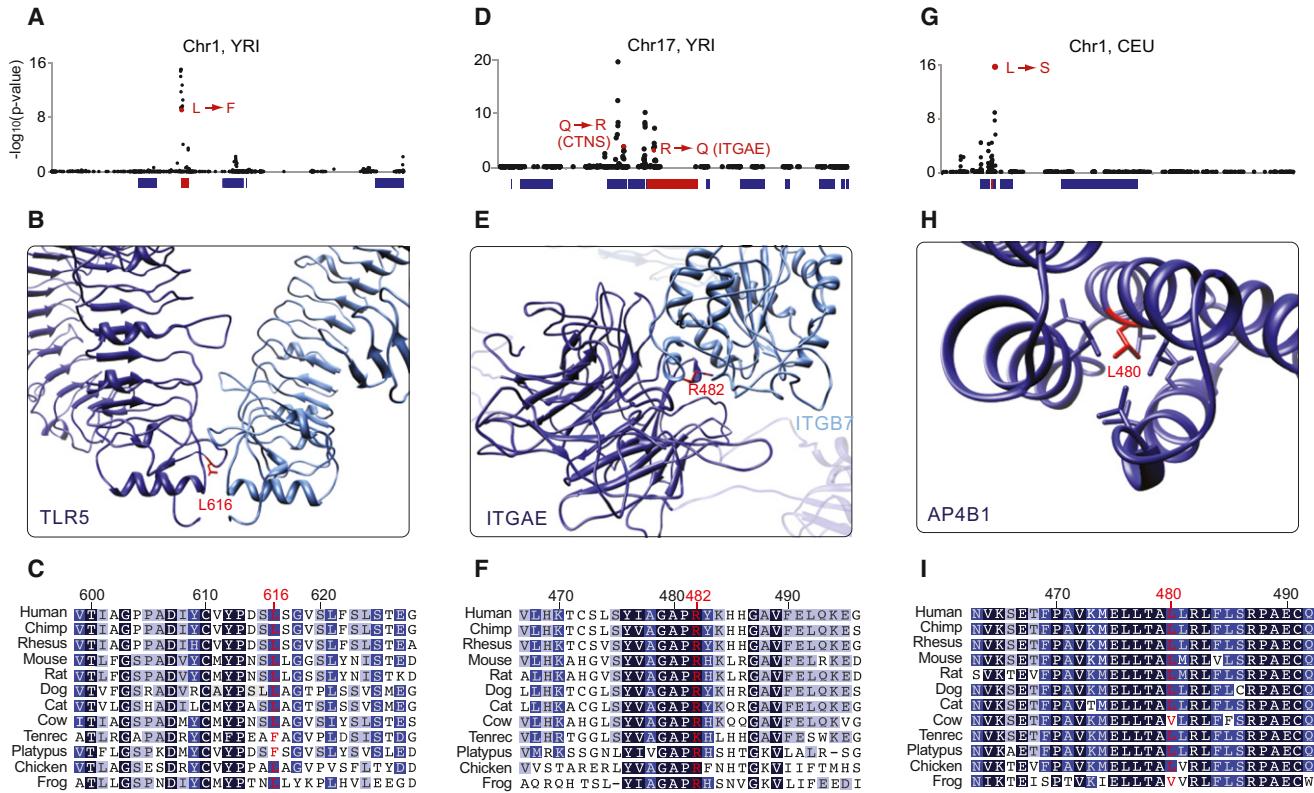
### 35 Candidate Adaptive Nonsynonymous Mutations

We examined the high-scoring variants and localized regions to identify both the biologically functional variants and the relevant pathways and phenotypes. We first focused on coding variation and identified 35 high-scoring nonsynonymous SNPs in 33 genes, of which seventeen are highly evolutionarily constrained (genomic evolutionary rate profiling [GERP] score greater than 2.0) (Cooper et al., 2005) (Figure 2 and Table S6). Of the 35, two have previously been characterized as adaptive mutations (in SCL24A5 and MATP, both associated with lighter skin pigmentation), and six have been associated with phenotypes in GWAS but not previously investigated as targets of selection. The latter include polymorphisms in EDAR (associated with hair thickness), ARHGEF3 (associated with greater bone mineral density), BTLA (rheumatoid arthritis), CTNS (cysteine metabolism defects), ITPR3 (type 1 diabetes and coronary aneurisms), and TLR5 (increased IFN $\gamma$  secretion). We performed further structural homology modeling and conservation analysis, which pinpointed possible functional roles for SNPs in several genes, including TLR5, ITGAE, and AP4B1.

It is striking that there are only 35 nonsynonymous variants in our entire list of candidates. Based on the genomic coverage of the 1000G data, we estimate that there are no more than 38 candidate causal nonsynonymous SNPs in the 412 candidate selected regions we analyzed (95% confidence interval, Experimental Procedures). These data suggest that only a minority of recent adaptations are due to amino-acid changes and that regulatory changes are likely to play a dominant role in recent human evolution.

### Numerous Candidate Adaptive Regulatory Elements

To begin our investigation of potential regulatory changes in recent human evolution, we compiled several published eQTL studies in the 1000G individuals (Montgomery et al., 2010; Pickrell et al., 2010; Stranger et al., 2007) and identified candidate variants in the 412 regions that have been associated with differences in gene expression (Figure 3A and Table S7). We identified 56 regions containing SNPs associated with expression levels of nearby genes in lymphoblast cells, a two-fold enrichment over the number expected by chance ( $p = 0.02$ ). In many of these cases, the top-scoring variants by CMS are the most strongly associated with expression levels. Fourteen high-scoring SNPs associated with expression differences lie in predicted transcription factor binding sites from the University of California Santa



**Figure 2. Candidate Nonsynonymous Mutations Identified by CMS**

(A–I) CMS identified high-scoring nonsynonymous mutations in the genes (A–C) *TLR5*, (D–F) *ITGAE*, and (G–I) *AP4B1*. (A, D, and G) CMS scores for all variants in the regions. High-scoring nonsynonymous variants and the genes in which they are located are presented in red. (B, E, and H) Homology modeling of the genes with the residue containing the candidate variants in red. (C, F, and I) The amino-acid sequence in 12 vertebrate species. The color of the residue indicates the conservation score (darker color indicates greater conservation). See also Table S6.

Cruz transcription factor binding site (UCSC TFBS) conserved track (Kent et al., 2002).

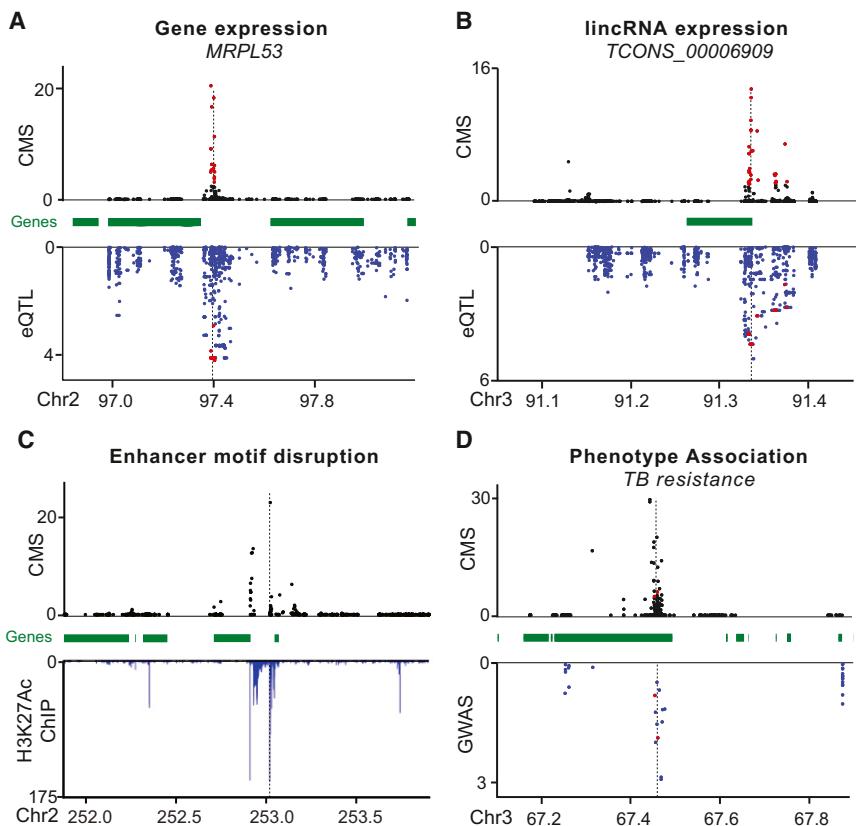
We extended our study of putatively selected regulatory variants to those affecting expression levels of lncRNAs, using recently published RNA sequencing data from 1000G samples (Experimental Procedures) (Montgomery et al., 2010; Pickrell et al., 2010). Three of the regions contain candidate variants under selection that are associated with differential expression of nearby lncRNAs (Figure 3B). To our knowledge, these three loci are the first candidates for recent selective pressure on mutations functionally affecting lncRNAs. LncRNAs themselves are often involved in orchestrating gene expression programs (Guttman et al., 2009; Guttman and Rinn, 2012).

Because available eQTL studies are somewhat underpowered and have only been carried out in lymphoblasts (likely missing many regulatory polymorphisms relevant in other cell types), we further used an alternative epigenetic-based strategy to identify potential regulatory variants. Using chromatin state predictions from (Ernst et al., 2011) we identified 335 SNPs (in 184 distinct candidate regions) that lie within predicted active enhancers or promoters and disrupt binding motifs of transcription factors known to be active in the cell type (Figure 3C and Table S8).

### Numerous Candidate Adaptive Variations Associated with Phenotypes

In parallel to our investigation of functional variants on a molecular level, we also characterized potential phenotypes and pathways linked to the candidate variants. As natural selection can only act on mutations that drive phenotypic variation, we examined polymorphisms that have been associated with a variety of traits. Using the National Human Genome Research Institute (NHGRI) Genome-Wide Association Study (GWAS) (Hindorff et al., 2009) database, we found 165 overlaps with trait-associated SNPs, including 11 regions that contain variants associated with height and pigmentation and 79 regions associated with infectious and autoimmune disease susceptibility (Table S9) (Davila et al., 2010; Fellay et al., 2007; Ge et al., 2009; Jallow et al., 2009; Kamatani et al., 2009; Mbarek et al., 2011; Png et al., 2011; Zhang et al., 2009).

We more closely examined one example of these GWAS overlaps: susceptibility to the mycobacterial pathogens *M. tuberculosis* and *M. leprae*, the causative agents of tuberculosis (TB) and leprosy, respectively. These pathogens have been a major source of morbidity and mortality and represent a possible selective pressure in human populations. In collaboration with the Wellcome Trust Case Control Consortium, we



**Figure 3. CMS Signals Overlapping Potential Regulatory Mutations, GWAS Signals, and Enhancers**

(A and B) Example of candidate selected variants associated with gene expression (A) and lincRNA expression (B). CMS scores (top) and eQTL p values (bottom) for all SNPs in the selected region. (C) Example of candidate selected variants that disrupt a putative enhancer element. CMS scores (top) and H3K27Ac ChIP-seq enrichment (bottom) from (Ernst et al., 2011).

(D) Example of candidate selected variants associated with TB resistance. CMS scores (top) and association test p values (bottom). Positions are given in centimorgans. High-scoring CMS variants that are with significant association scores are shown in red.

See also Tables S7, S8, and S9.

analyzed a TB genome-wide association study (GWAS) conducted in West Africa and a leprosy GWAS conducted in northern India to find overlap between the localized selected regions and loci associated with TB and leprosy susceptibility (Thye et al., 2010; Wong et al., 2010) (see Experimental Procedures). Five variants that are associated with resistance to these infectious diseases at  $p < 10^{-5}$  fall within the fine-mapped CMS-identified regions (Figure 3D and Table S4). The locus with the strongest association with resistance to leprosy ( $p = 1.25 \times 10^{-6}$ ) contains *SLC24A5*, also known to influence skin pigmentation. Other loci suggested to be associated with leprosy include *ALMS1*, and with tuberculosis resistance include *CCR9*, *CXCR4*, and *VDR*, all of which play a role in immune response or pathogen binding (Liu et al., 2004).

#### A Candidate Adaptive Mutation in TLR5 Leads to Diminished NF- $\kappa$ B Signaling

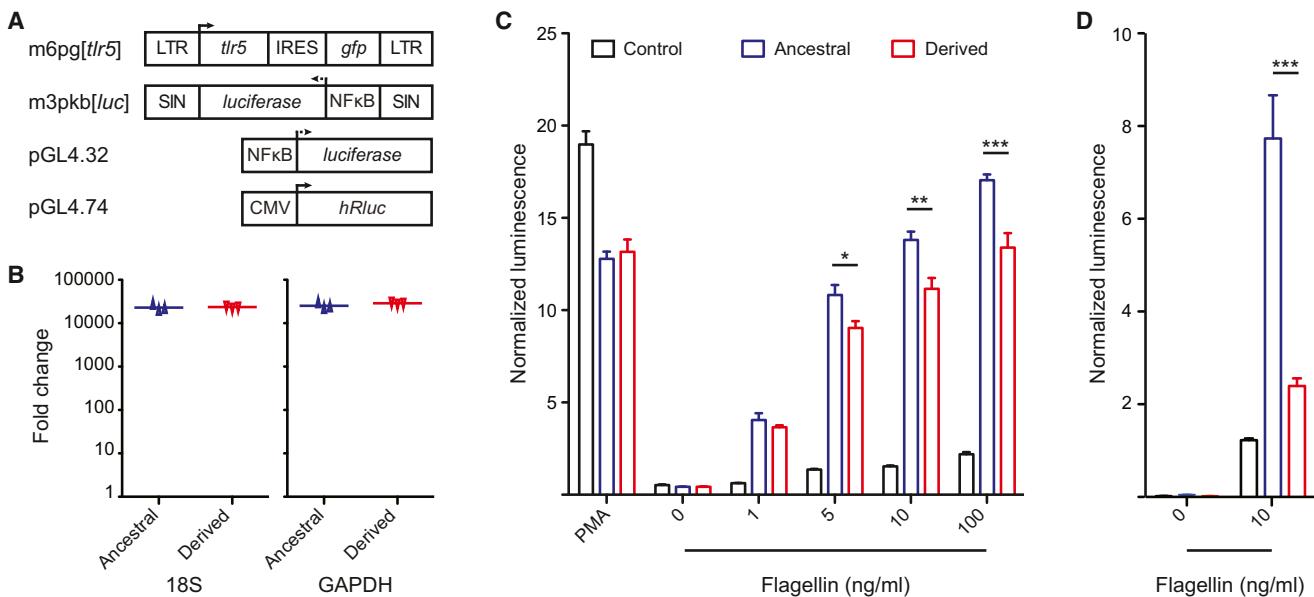
We chose to functionally characterize one of the most promising nonsynonymous candidates of selection, a leucine-to-phenylalanine variant (L616F) in *TLR5*, detected in the YRI. This is the only coding SNP in the *TLR5* region with evidence of selection, and none of the high-scoring SNPs outside the open reading frame (ORF) appear to affect Toll-like receptor 5 (TLR5) expression or disrupt regulatory regions. TLR5 is a well-described Toll-like receptor that plays a crucial role in the immunological clearance of bacterial pathogens. Its ligand, flagellin, is the principal component of the bacterial flagellum

and is one of the most abundantly expressed proteins in nearly all flagellated bacteria. Receptor-ligand interaction of TLR5 activates cells of the immune system via the NF- $\kappa$ B pathway, leading to the production of various proinflammatory mediators. Several polymorphisms in *TLR5* have been associated with differential responses to infectious diseases including Legionnaire's disease (Hawn et al., 2003) and neonatal sepsis (Abu-Maziad et al., 2010).

Using structural modeling, we found that the TLR5 L616F variant is predicted to be located in the conserved ectodomain responsible for dimerization and activation of the receptor (Figures 2B and 2C). To test the cellular effect of the TLR5 L616F mutation, we created stable cell lines carrying either the ancestral (L) or the derived (F) form of TLR5 in two different cell types and measured NF- $\kappa$ B activation in the presence of increasing amounts of flagellin (Experimental Procedures and Figure 4). In both cell types, we found that the derived TLR5-616F allele produced significantly reduced NF- $\kappa$ B signaling in response to flagellin relative to the ancestral TLR5-616L allele (Figures 4C and 4D).

#### DISCUSSION

The promise of the genomic age for elucidating human evolution has not yet been realized, in part due to the large size of regions identified as targets of selection, each of which can contain thousands of candidate causal variants, and in part due to the incompleteness of genotype data. Drawing on full-genome sequence data from 1000G and on the CMS method, this paper presents a comprehensive catalog of potential human adaptive mutations, instead of genomic regions. Each fine-mapped region contains 20–100 candidate variants, a small enough number to be tractable for functional characterization. As causal variants under selection typically have 10–50 perfect proxy variants, we are already near the limit of the power of population genetic tests



**Figure 4. The Derived Form of TLR5 Leads to a Diminished NF-κB Response**

(A) Structure of retroviral vectors containing ancestral and derived TLR5 alleles and plasmid reporter construct transduced into 293FT and Jurkat cells. (B) Expression of TLR5 relative to nontransduced cells in transduced 293FT cells, given as normalized levels to either the 18S ribosomal subunit or GAPDH. (C) NF-κB reporter activity of 293FT cells transduced with ancestral or derived TLR5 allele or empty vector control 24 hr after stimulation with varying amounts of flagellin, normalized against the renilla luciferase signal. (D) NF-κB reporter activity of Jurkat cells transduced with ancestral or derived TLR5 allele or empty vector control 24 hr after stimulation with flagellin, normalized to nonspecific activation with PMA/ionomycin. Control lane represents cells transduced with empty m6pg vector. Error bars represent the SEM of at least three independent experiments.

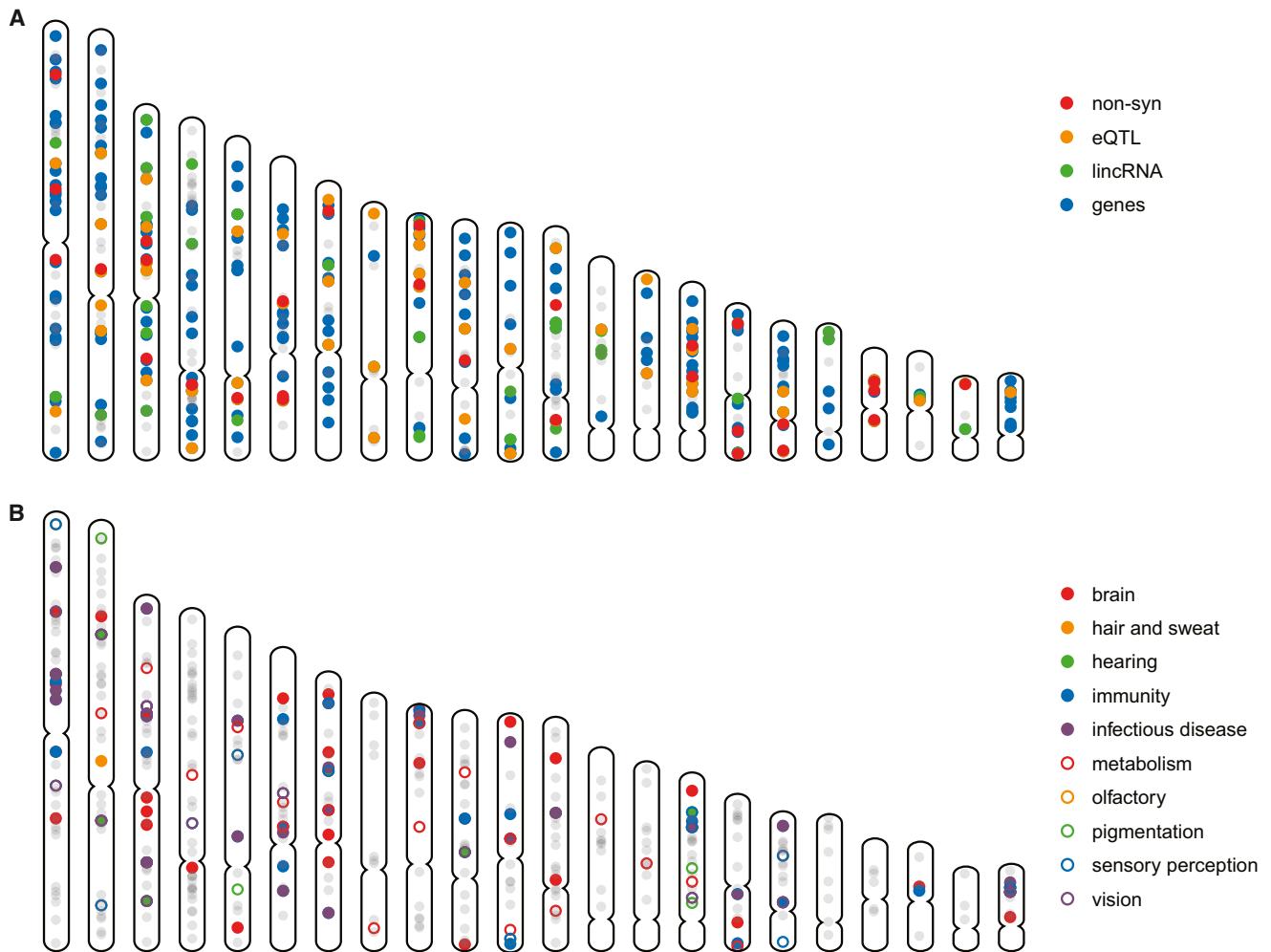
to pinpoint the causal variant (Grossman et al., 2010). We computationally annotated all candidates and provide a proof-of-principle example of functional validation, creating a rich resource for future studies of human adaptations.

Many of the variants thus identified are associated with pathways that have emerged as targets of the strongest selective pressures on humans in recent history; the relevant traits include skin color, metabolism, and infectious disease resistance (Figure 5). In addition to the phenotypic associations and gene enrichment in these pathways discussed above, several of the eQTL SNPs regulate the expression of genes in these pathways, such as IVD, ACAS2, and CTNS (involved in metabolism) and BLK (involved in immune function). Many mutations fall in and around genes encoding the receptors or enzymes that modify the receptors for some of the most devastating pathogens in human history, including *RHOA* and *OTUB1* (*Yersinia pestis*) (Edelmann et al., 2010) and *DAG1* (*Mycobacterium leprae*), *TLR5* (*Salmonella typhimurium* and others), *LARGE* (Lassa virus) (Kunz et al., 2005), *DARC* (*Plasmodium vivax* malaria) (Sabeti et al., 2006), *PVRL4* (measles virus), *VDR* (*Mycobacterium tuberculosis*), *TPST1* (HIV) (Farzan et al., 1999), and *CXCR4* (HIV). New pathways under selection are also coming to light: for example, in this issue of Cell Kamberov et al. (2013) elucidates selection on a nonsynonymous mutation in *EDAR*, which leads to a number of pleiotropic traits including altered hair and sweat gland formation.

Our data support the mounting evidence that a great deal of recent human adaptation and phenotypic variation is based in

regulatory regions (Hindorff et al., 2009; Lindblad-Toh et al., 2011; Vernot et al., 2012; Wang et al., 1995). Less than 10% of our fine-mapped regions contained high-scoring nonsynonymous SNPs; candidate selected SNPs are enriched for eQTLs and include many mutations that disrupt transcription factor motifs in enhancers and promoters. Motifs for transcription factors involved in a number of different processes are disrupted, including STATs, Jun, GATAs, C/EBP, PPAR $\gamma$ , ETS, and IRFs. In several cases, the motif for a cell-specific transcription factor is disrupted in a cell-specific enhancer, for example an LXR:RXR motif in a hepatocyte-specific enhancer or a PU.1 motif in a monocyte-specific DNase HS site. The magnitude of the change in binding affinity varied from a minimum change in LOD score of 0.3 to a maximum of 12. More complete characterization of the regulatory variants using high-throughput cellular assays and eQTL studies in additional individuals may be illuminating.

Given the bounds on population genetic approaches to fine-map signals of selection and the limitations of current functional annotations, the true adaptive mutation must ultimately be distinguished using functional approaches. This is a challenge, especially for regions identified through genome-wide scans instead of based on a prior hypothesis of an adaptive pressure (e.g., malaria and lactose tolerance). It is impracticable to assay each variant in every possible cell type and process, and furthermore, even functional variants need not be causal. Although there is no way to prove what evolution did, even if we could go back in time to observe it, the standard in the evolutionary



**Figure 5. Characterization of Candidate Regions and Variants**

(A and B) All candidate regions in the genome are shown in gray. (A) Candidate functional elements in localized regions, including regions with genes (blue), eQTLs (orange), long noncoding RNAs (green), and nonsynonymous variants (red). (B) Regions with genes relating to potential selective pressures, such as metabolism (red circle), infectious disease (purple), brain development (red), hearing (green), and hair and sweat (orange).

See also Tables S2, S3, S4, and S5.

genomics field for establishing a mutation as having caused selection is strong statistical evidence of selection plus a phenotypic effect likely to enhance survival.

We chose one of the candidate variants, a nonsynonymous mutation in *TLR5*, to characterize experimentally. The derived allele in *TLR5* with evidence of selection leads to diminished NF- $\kappa$ B signaling during bacterial infections. Intriguingly, another allele that decreases the function of *TLR5* (a nonsense variant, *TLR5-392STOP*) has previously been reported to reach a frequency of 10% in European populations (Barreiro et al., 2009). The existence of common variants that decrease *TLR5* signaling suggests that modulating *TLR5* signaling may be advantageous in certain environments. Indeed, decreasing NF- $\kappa$ B signaling can have a protective effect in several bacterial infections, most significantly in bacterial sepsis (Koedel et al., 2000; Okugawa et al., 2006). Furthermore, the pathogen *Salmonella typhimurium* requires activated lamina propria cells (LPCs)

in the intestinal epithelium to invade a host and is consequently unable to infect mice with deficient TLR5 signaling (Uematsu et al., 2006). In a human population constantly exposed to high levels of bacterial antigens, a TLR5 variant with reduced NF- $\kappa$ B activation may well confer a fitness advantage.

An accompanying article from Kamberov et al. (2013) models an adaptive human variant of *EDAR* in mice and characterizes its phenotype and evolutionary origins in humans. *EDARV370A*, one of the 35 nonsynonymous variants detected by CMS in 1000G data, likely emerged in central China ~30,000 years ago and leads to increased sweat gland number and scalp hair thickness in mice and humans. *TLR5L616F* and *EDARV370A* demonstrate the power of our framework to move from genomic scans to the characterization of a novel adaptive mutation and elucidation of distinct mechanisms of evolution.

This paper, in conjunction with the accompanying paper on *EDAR*, represents a decisive shift for the field of evolutionary

genomics, moving from hypothesis-driven to hypothesis-generating science. We further provide a comprehensive list of candidate adaptive *mutations* driving recent human selective sweeps that lay the foundation for myriad future functional studies. The data from the 1000G Project, along with functional annotations, are available on a genome-wide browser, together with software to compute CMS on any data set (<http://www.broadinstitute.org/mpg/cms>). In the years ahead, unprecedented data availability and collaborations across multiple disciplines from molecular, developmental, and computational biology to history and anthropology, promise to bring key recent events that have shaped our species to light.

## EXPERIMENTAL PROCEDURES

### Simulations

We used the simulations described earlier in Grossman et al. (2010), with one change: a coding error was fixed in the code that simulated gene conversion during a selective sweep (neutral simulations were unaffected).

### CMS

Two versions of the CMS test were used: the original (within-region) CMS test (Grossman et al., 2010) for localizing the selected variant within a candidate region, and a modified test (denoted CMS<sub>GW</sub>) for identifying candidate regions within the genome.

When using CMS to localize regions, we used the distribution of neutral SNPs within 500 kb of selected SNPs as the “unselected” distribution and assumed exactly one selected SNP per region. To use CMS as a genome-wide method to detect selected regions, we made the following modifications:

- (1) SNPs in neutral regions were used as the “unselected” distribution
- (2) We did not assume any prior hypothesis about how many SNPs are under selection. Therefore instead of calculating the posterior probability, we calculated the Bayes factor for each test

$$BF_t = \frac{P(v_t \in bin_{t,k} | selected)}{P(v_t \in bin_{t,k} | unselected)}$$

and defined the composite score as the product of the Bayes factor of each test:

$$CMS_{GW} = \prod_{t \in tests} BF_t$$

- (3) Scores were normalized to neutral simulations (for simulated data) or to the whole genome (for real data), rather than within each region.
- (4) Bin boundaries were adjusted as described earlier.

We identified 100 kb regions in which 30% of SNPs had a normalized score above 3, a threshold that corresponded to a 0.1% false positive rate (FPR) in simulations (i.e., in 1,000 neutral simulations of a 1 Mb region no more than 1 contained a 100 kb region meeting this criterion; the upper bound of the 95% binomial confidence interval for the FPR is 0.6%), and used this threshold to detect selected regions in the 1000G data. We note that because the 1000G data include 2.42 Gb, we expect 24 false positive regions at this threshold.

The code used for CMS analysis is available at <http://www.broadinstitute.org/mpg/cms>.

### 1000G

Quality-controlled phased SNP and indel calls for the CEU, YRI and CHB+JPT populations released by the low-coverage portion of the 1000G project were downloaded from [ftp://ftp.1000genomes.ebi.ac.uk/vol1/ftp/pilot\\_data/release/2010\\_03/pilot1/](ftp://ftp.1000genomes.ebi.ac.uk/vol1/ftp/pilot_data/release/2010_03/pilot1/), representing the March 2010 data release. All genetic variants with more than two alleles were converted to biallelic variants, by mapping all alleles to two alleles while preserving alleles’ ancestral state where

known. Ancestral state was taken from the ancestral state data released by the 1000G Project at [ftp://ftp.1000genomes.ebi.ac.uk/vol1/ftp/pilot\\_data/technical/reference/ancestral\\_alignments](ftp://ftp.1000genomes.ebi.ac.uk/vol1/ftp/pilot_data/technical/reference/ancestral_alignments), constructed from a four-way alignment of human, chimp, orangutan, and rhesus macaque. Monomorphic SNPs omitted from 1000G data but present in HapMap Phase II data were added back into the data.

### Gene Pathway Enrichment

To test for enrichment of different classes of functional variants, we picked random sets of 412 nonoverlapping regions in the genome that matched our selected regions in size and compared the number overlaps in our selected regions to the number in the random regions.

We also manually defined functional categories that previous literature suggests may have played an important role in recent human adaptation (e.g., skin pigmentation and immune system processes). We then used INRICHv1.0 to test for enrichment of these pathways. INRICH uses a two-step permutation algorithm to test for enrichment of pathways defined by the user or derived from published databases and corrects the number of gene sets tested, the size of each geneset, and the number of SNPs and genes within our selected regions. See <http://atgu.mgh.harvard.edu/inrich/> for more information.

### Protein Structure Modeling

We used Modeller9v8 to generate homology models of proteins in which we found a high-scoring nonsynonymous SNP. Sequences similar to the target sequence were selected as templates for homology modeling, and the optimum model was selected as the one with the lowest energy (DOPE) score. For TLR5, we used a published computationally derived model of human TLR5, provided by Wei and colleagues (Wei et al., 2011).

### Cell Lines

Transgenes carrying either the ancestral (*tlr5a*; leucine) or derived (*tlr5d*; phenylalanine) form of TLR5 were synthesized and cloned into the retroviral vector m6pg carrying GFP as a transgene and transduced into 293FT and Jurkat cell to create stable cell lines. TLR5 expression was measured by qPCR.

### 293FT and Jurkat Luciferase Assays

NF-κB activity in 293FT cells was measured using pGL4.32 and pGL4.74 (Promega) and in Jurkat cells using the retroviral reporter m3pkb[*luc*] carrying an NF-κB inducible luciferase reporter (Loizou et al., 2011).

293FT cells expressing either the ancestral or derived forms of TLR5 were transfected with pGL4.32 and pGL4.74 (Promega). Twenty-six hours after transfection, the cells were stimulated for an additional 24 hr with 800 ng/ml PMA or increasing levels of flagellin at 1, 5, 10, or 100 ng/ml, and luciferase and renilla luminescence was measured in a Top Count machine.

293FT cells were transfected with 6.6 μg each of m3p[*luc*], retroviral packaging pCL-Eco, and viral envelope Vsv-g DNA. Cells were incubated for 20 hr and supernatant containing viral particles was collected.

The m3pkb[*luc*] viral supernatant and protamine sulfate were added to Jurkat cells stably expressing either the ancestral or derived forms of TLR5 and spun at 400 × g for 2 hr at 32°C. Plates were incubated for 26 hr (replacing media after 20 hr). Cells were stimulated for 24 hr with 400 ng/ml PMA and 1.5 μg/ml ionomycin or 10 ng/ml flagellin and firefly luminescence was measured.

### Gene and lincRNA eQTL Analysis

We obtained expression intensities of 47,293 probes representing the majority of human genes from Stranger et al. (Stranger et al., 2007) and downloaded the normalized gene expression levels for 22,032 genes in the YRI individuals measured by RNA seq by Pickrell et al. (2010).

To investigate lincRNA expression, RNA reads for YRI (Pickrell et al., 2010) were obtained from [http://eqtl.uchicago.edu/RNA\\_Seq\\_data/](http://eqtl.uchicago.edu/RNA_Seq_data/) and for CEU (Montgomery et al., 2010) from [http://jungle.unige.ch/rnaseq\\_CEU60/](http://jungle.unige.ch/rnaseq_CEU60/) and aligned to hg19 using BWA (Li and Durbin, 2009). Aligned reads were counted across the 4,421 previously detected regions of interest. Reads were RPKM normalized against both the length of the region and the total read count in the lane (Mortazavi et al., 2008) to provide a baseline expression

level for each region. LncRNAs with nonzero expression in at least half of the individuals in a population were analyzed. All nonzero expression levels were quantile-normalized within each population in order to produce a normal distribution of expression.

Normalized read counts from each gene or lincRNA were tested as quantitative traits in a standard association test with SNPs 1 Mb. SNPs with significant association p values were overlapped with the selected regions.

### TB and Leprosy Association Studies

TB susceptibility data were obtained from the Wellcome Trust Case Control Consortium study in the Gambia (Thye et al., 2010) with 1,498 confirmed TB cases and 1,496 controls, genotyped on the Affymetrix GeneChip 500K Array comprising 500,568 SNPs using the CHIAMO algorithm. The primary analysis focused on single-locus tests of association using 1,320 TB cases compared to 1,384 Gambian controls for all 405,226 SNPs passing QC filters with a study-wide MAF > 1%. The trend test was performed in a logistic regression modeling framework, which was adjusted for three axes of multidimensional scaling, by inclusion as covariates in the logistic regression model, reducing the overdispersion of trend tests from  $\lambda = 1.13$  (no adjustment) to  $\lambda = 1.05$ .

Leprosy susceptibility data were obtained from the host genetics study of leprosy in Indians (Wong et al., 2010), consisting of 258 confirmed cases of leprosy and 300 controls from New Delhi. All individuals in this study were genotyped with the Illumina IBC gene-centric 50K array. Multidimensional scaling (MDS) and principal component analysis (PCA) were carried out with PLINK and EIGENSTRAT to remove population outliers. A total of 209 leprosy cases and 239 controls were carried forward for analysis after quality control filters. The primary test of association in the New Delhi and Kolkata cohorts was carried out with the Pearson's  $\chi^2$  allelic test, Cochran-Armitage trend test and logistic regression.

### SUPPLEMENTAL INFORMATION

Supplemental Information includes Extended Experimental Procedures, one figure, and nine tables and can be found with this article online at <http://dx.doi.org/10.1016/j.cell.2013.01.035>.

### ACKNOWLEDGMENTS

P.C.S. and her lab are supported by a Burroughs Wellcome Fund Career Award, a Packard Foundation Fellowship in Science and Engineering, a Broad Institute SPARC award, an NIH Innovator Award 1DP2OD006514-01, and BAA-NIAID-DAIT-NIHA12009061. S.R.G. is supported by NIGMS T32GM007753, K.G.A. by a Carlsberg Foundation fellowship, D.J.P. by NSF, and E.K.K. by an American Cancer Society Fellowship. The mycobacterial disease studies analyzed were supported by funding from the Wellcome Trust, the UK Medical Research Council, the UK National Institute for Health Research, and the European Commission; we thank the many collaborators who contributed to generating these data sets. We would like to thank S. Hart for help with figures, C. Edwards for reviews of the text, and L. Ward and C. O'Dushlaine for technical guidance.

Received: October 6, 2012

Revised: January 15, 2013

Accepted: January 22, 2013

Published: February 14, 2013

### REFERENCES

- 1000 Genomes Project Consortium, Abecasis, G.R., Altshuler, D., Auton, A., Brooks, L.D., Durbin, R.M., Gibbs, R.A., Hurles, M.E., and McVean, G.A. (2010). A map of human genome variation from population-scale sequencing. *Nature* 467, 1061–1073.
- Abu-Maziad, A., Schaa, K., Bell, E.F., Dagle, J.M., Cooper, M., Marazita, M.L., and Murray, J.C. (2010). Role of polymorphic variants as genetic modulators of infection in neonatal sepsis. *Pediatr. Res.* 68, 323–329.
- Akey, J.M. (2009). Constructing genomic maps of positive selection in humans: where do we go from here? *Genome Res.* 19, 711–722.
- Akey, J.M., Zhang, G., Zhang, K., Jin, L., and Shriver, M.D. (2002). Interrogating a high-density SNP map for signatures of natural selection. *Genome Res.* 12, 1805–1814.
- Barreiro, L.B., Ben-Ali, M., Quach, H., Laval, G., Patin, E., Pickrell, J.K., Boucher, C., Tchitch, M., Neyrolles, O., Gicquel, B., et al. (2009). Evolutionary dynamics of human Toll-like receptors and their different contributions to host defense. *PLoS Genet.* 5, e1000562.
- Bersaglieri, T., Sabeti, P.C., Patterson, N., Vanderploeg, T., Schaffner, S.F., Drake, J.A., Rhodes, M., Reich, D.E., and Hirschhorn, J.N. (2004). Genetic signatures of strong recent positive selection at the lactase gene. *Am. J. Hum. Genet.* 74, 1111–1120.
- Bustamante, C.D., Fledel-Alon, A., Williamson, S., Nielsen, R., Hubisz, M.T., Golanowski, S., Tanenbaum, D.M., White, T.J., Sninsky, J.J., Hernandez, R.D., et al. (2005). Natural selection on protein-coding genes in the human genome. *Nature* 437, 1153–1157.
- Cabili, M.N., Trapnell, C., Goff, L., Koziol, M., Tazon-Vega, B., Regev, A., and Rinn, J.L. (2011). Integrative annotation of human large intergenic noncoding RNAs reveals global properties and specific subclasses. *Genes Dev.* 25, 1915–1927.
- Cooper, G.M., Stone, E.A., Asimenos, G., Green, E.D., Batzoglou, S., and Sidow, A.; NISC Comparative Sequencing Program. (2005). Distribution and intensity of constraint in mammalian genomic sequence. *Genome Res.* 15, 901–913.
- Currat, M., Trabuchet, G., Rees, D., Perrin, P., Harding, R.M., Clegg, J.B., Langaney, A., and Excoffier, L. (2002). Molecular analysis of the beta-globin gene cluster in the Niokholo Mandenka population reveals a recent origin of the beta(S) Senegal mutation. *Am. J. Hum. Genet.* 70, 207–223.
- Davila, S., Wright, V.J., Khor, C.C., Sim, K.S., Binder, A., Breunis, W.B., Inwald, D., Nadel, S., Betts, H., Carroll, E.D., et al.; International Meningococcal Genetics Consortium. (2010). Genome-wide association study identifies variants in the CFH region associated with host susceptibility to meningococcal disease. *Nat. Genet.* 42, 772–776.
- Edelmann, M.J., Kramer, H.B., Altun, M., and Kessler, B.M. (2010). Post-translational modification of the deubiquitinating enzyme otubain 1 modulates active RhoA levels and susceptibility to Yersinia invasion. *FEBS J.* 277, 2515–2530.
- Ernst, J., Kheradpour, P., Mikkelsen, T.S., Shoresh, N., Ward, L.D., Epstein, C.B., Zhang, X., Wang, L., Issner, R., Coyne, M., et al. (2011). Mapping and analysis of chromatin state dynamics in nine human cell types. *Nature* 473, 43–49.
- Farzan, M., Mirzabekov, T., Kolchinsky, P., Wyatt, R., Cayabyab, M., Gerard, N.P., Gerard, C., Sodroski, J., and Choe, H. (1999). Tyrosine sulfation of the amino terminus of CCR5 facilitates HIV-1 entry. *Cell* 96, 667–676.
- Fellay, J., Shianna, K.V., Ge, D., Colombo, S., Ledergerber, B., Weale, M., Zhang, K., Gumbi, C., Castagna, A., Cossarizza, A., et al. (2007). A whole-genome association study of major determinants for host control of HIV-1. *Science* 317, 944–947.
- Frazer, K.A., Ballinger, D.G., Cox, D.R., Hinds, D.A., Stuve, L.L., Gibbs, R.A., Belmont, J.W., Boudreau, A., Hardenbol, P., Leal, S.M., et al.; International HapMap Consortium. (2007). A second generation human haplotype map of over 3.1 million SNPs. *Nature* 449, 851–861.
- Ge, D., Fellay, J., Thompson, A.J., Simon, J.S., Shianna, K.V., Urban, T.J., Heinzen, E.L., Qiu, P., Bertelsen, A.H., Muir, A.J., et al. (2009). Genetic variation in IL28B predicts hepatitis C treatment-induced viral clearance. *Nature* 461, 399–401.
- Grossman, S.R., Shlyakhter, I., Karlsson, E.K., Byrne, E.H., Morales, S., Frieden, G., Hostetter, E., Angelino, E., Garber, M., Zuk, O., et al. (2010). A composite of multiple signals distinguishes causal variants in regions of positive selection. *Science* 327, 883–886.

- Gusev, A., Palamara, P.F., Aponte, G., Zhuang, Z., Darvasi, A., Gregersen, P., and Pe'er, I. (2012). The architecture of long-range haplotypes shared within and across populations. *Mol. Biol. Evol.* 29, 473–486.
- Guttman, M., and Rinn, J.L. (2012). Modular regulatory principles of large non-coding RNAs. *Nature* 482, 339–346.
- Guttman, M., Amit, I., Garber, M., French, C., Lin, M.F., Feldser, D., Huarte, M., Zuk, O., Carey, B.W., Cassady, J.P., et al. (2009). Chromatin signature reveals over a thousand highly conserved large non-coding RNAs in mammals. *Nature* 458, 223–227.
- Hawn, T.R., Verbon, A., Lettinga, K.D., Zhao, L.P., Li, S.S., Laws, R.J., Skerrett, S.J., Beutler, B., Schroeder, L., Nachman, A., et al. (2003). A common dominant TLR5 stop codon polymorphism abolishes flagellin signaling and is associated with susceptibility to legionnaires' disease. *J. Exp. Med.* 198, 1563–1572.
- Hernandez, R.D., Kelley, J.L., Elyashiv, E., Melton, S.C., Auton, A., McVean, G., Sella, G., and Przeworski, M.; 1000 Genomes Project. (2011). Classic selective sweeps were rare in recent human evolution. *Science* 331, 920–924.
- Hindorff, L.A., Sethupathy, P., Junkins, H.A., Ramos, E.M., Mehta, J.P., Collins, F.S., and Manolio, T.A. (2009). Potential etiologic and functional implications of genome-wide association loci for human diseases and traits. *Proc. Natl. Acad. Sci. USA* 106, 9362–9367.
- Jallow, M., Teo, Y.Y., Small, K.S., Rockett, K.A., Deloukas, P., Clark, T.G., Kivinen, K., Bojang, K.A., Conway, D.J., Pinder, M., et al.; Wellcome Trust Case Control Consortium; Malaria Genomic Epidemiology Network. (2009). Genome-wide and fine-resolution association analysis of malaria in West Africa. *Nat. Genet.* 41, 657–665.
- Kamatani, Y., Wattanapokayakit, S., Ochi, H., Kawaguchi, T., Takahashi, A., Hosono, N., Kubo, M., Tsunoda, T., Kamatani, N., Kurnada, H., et al. (2009). A genome-wide association study identifies variants in the HLA-DP locus associated with chronic hepatitis B in Asians. *Nat. Genet.* 41, 591–595.
- Kamberov, Y.G., Wang, S., Tan, J., Gerbault, P., Wark, A., Tan, L., Yang, Y., Li, S., Tang, K., Chen, H., et al. (2013). Modeling Recent Human Evolution in Mice by Expression of a Selected EDAR Variant. *Cell* 152, this issue, 691–702.
- Kent, W.J., Sugnet, C.W., Furey, T.S., Roskin, K.M., Pringle, T.H., Zahler, A.M., and Haussler, D. (2002). The human genome browser at UCSC. *Genome Res.* 12, 996–1006.
- Koedel, U., Bayerlein, I., Paul, R., Sporer, B., and Pfister, H.W. (2000). Pharmacologic interference with NF- $\kappa$ B activation attenuates central nervous system complications in experimental Pneumococcal meningitis. *J. Infect. Dis.* 182, 1437–1445.
- Kunz, S., Rojek, J.M., Kanagawa, M., Spiropoulou, C.F., Barresi, R., Campbell, K.P., and Oldstone, M.B. (2005). Posttranslational modification of alpha-dystroglycan, the cellular receptor for arenaviruses, by the glycosyltransferase LARGE is critical for virus binding. *J. Virol.* 79, 14282–14296.
- Lamason, R.L., Mohideen, M.A., Mest, J.R., Wong, A.C., Norton, H.L., Aros, M.C., Juryne, M.J., Mao, X., Humphreville, V.R., Humbert, J.E., et al. (2005). SLC24A5, a putative cation exchanger, affects pigmentation in zebrafish and humans. *Science* 310, 1782–1786.
- Li, H., and Durbin, R. (2009). Fast and accurate short read alignment with Burrows-Wheeler transform. *Bioinformatics* 25, 1754–1760.
- Lindblad-Toh, K., Garber, M., Zuk, O., Lin, M.F., Parker, B.J., Washietl, S., Kheradpour, P., Ernst, J., Jordan, G., Mauceli, E., et al.; Broad Institute Sequencing Platform and Whole Genome Assembly Team Baylor College of Medicine Human Genome Sequencing Center Sequencing Team Genome Institute at Washington University. (2011). A high-resolution map of human evolutionary constraint using 29 mammals. *Nature* 478, 476–482.
- Liu, W., Cao, W.C., Zhang, C.Y., Tian, L., Wu, X.M., Habbema, J.D., Zhao, Q.M., Zhang, P.H., Xin, Z.T., Li, C.Z., and Yang, H. (2004). VDR and NRAMP1 gene polymorphisms in susceptibility to pulmonary tuberculosis among the Chinese Han population: a case-control study. *Int. J. Tuberc. Lung Dis.* 8, 428–434.
- Loizou, L., Andersen, K.G., and Betz, A.G. (2011). Foxp3 interacts with c-Rel to mediate NF- $\kappa$ B repression. *PLoS ONE* 6, e18670.
- Mbarek, H., Ochi, H., Urabe, Y., Kumar, V., Kubo, M., Hosono, N., Takahashi, A., Kamatani, Y., Miki, D., Abe, H., et al. (2011). A genome-wide association study of chronic hepatitis B identified novel risk locus in a Japanese population. *Hum. Mol. Genet.* 20, 3884–3892.
- Montgomery, S.B., Sammeth, M., Gutierrez-Arcelus, M., Lach, R.P., Ingle, C., Nisbett, J., Guigo, R., and Dermitzakis, E.T. (2010). Transcriptome genetics using second generation sequencing in a Caucasian population. *Nature* 464, 773–777.
- Mortazavi, A., Williams, B.A., McCue, K., Schaeffer, L., and Wold, B. (2008). Mapping and quantifying mammalian transcriptomes by RNA-Seq. *Nat. Methods* 5, 621–628.
- Ohashi, J., Naka, I., Patarapotkul, J., Hananantachai, H., Brittenham, G., Looareesuwan, S., Clark, A.G., and Tokunaga, K. (2004). Extended linkage disequilibrium surrounding the hemoglobin E variant due to malarial selection. *Am. J. Hum. Genet.* 74, 1198–1208.
- Okugawa, S., Yanagimoto, S., Tsukada, K., Kitazawa, T., Koike, K., Kimura, S., Nagase, H., Hirai, K., and Ota, Y. (2006). Bacterial flagellin inhibits T cell receptor-mediated activation of T cells by inducing suppressor of cytokine signalling-1 (SOCS-1). *Cell. Microbiol.* 8, 1571–1580.
- Pickrell, J.K., Coop, G., Novembre, J., Kudaravalli, S., Li, J.Z., Absher, D., Srivivasan, B.S., Barsh, G.S., Myers, R.M., Feldman, M.W., and Pritchard, J.K. (2009). Signals of recent positive selection in a worldwide sample of human populations. *Genome Res.* 19, 826–837.
- Pickrell, J.K., Marioni, J.C., Pai, A.A., Degner, J.F., Engelhardt, B.E., Nkadori, E., Veyrieras, J.B., Stephens, M., Gilad, Y., and Pritchard, J.K. (2010). Understanding mechanisms underlying human gene expression variation with RNA sequencing. *Nature* 464, 768–772.
- Png, E., Thalamuthu, A., Ong, R.T., Snippe, H., Boland, G.J., and Seielstad, M. (2011). A genome-wide association study of hepatitis B vaccine response in an Indonesian population reveals multiple independent risk variants in the HLA region. *Hum. Mol. Genet.* 20, 3893–3898.
- Sabeti, P.C., Schaffner, S.F., Fry, B., Lohmueller, J., Varilly, P., Shamovsky, O., Palma, A., Mikkelsen, T.S., Altshuler, D., and Lander, E.S. (2006). Positive natural selection in the human lineage. *Science* 312, 1614–1620.
- Sabeti, P.C., Varilly, P., Fry, B., Lohmueller, J., Hostetter, E., Cotsapas, C., Xie, X., Byrne, E.H., McCarroll, S.A., Gaudet, R., et al.; International HapMap Consortium. (2007). Genome-wide detection and characterization of positive selection in human populations. *Nature* 449, 913–918.
- Stranger, B.E., Nica, A.C., Forrest, M.S., Dimas, A., Bird, C.P., Beazley, C., Ingle, C.E., Dunning, M., Flücke, P., Koller, D., et al. (2007). Population genomics of human gene expression. *Nat. Genet.* 39, 1217–1224.
- Thye, T., Vannberg, F.O., Wong, S.H., Owusu-Dabo, E., Osei, I., Gyapong, J., Sirugo, G., Sisay-Joof, F., Enimil, A., Chinbuah, M.A., et al.; African TB Genetics Consortium Wellcome Trust Case Control Consortium. (2010). Genome-wide association analyses identifies a susceptibility locus for tuberculosis on chromosome 18q11.2. *Nat. Genet.* 42, 739–741.
- Uematsu, S., Jang, M.H., Chevrier, N., Guo, Z., Kumagai, Y., Yamamoto, M., Kato, H., Sougawa, N., Matsui, H., Kuwata, H., et al. (2006). Detection of pathogenic intestinal bacteria by Toll-like receptor 5 on intestinal CD11c+ lamina propria cells. *Nat. Immunol.* 7, 868–874.
- Vernot, B., Stergachis, A.B., Maurano, M.T., Vierstra, J., Neph, S., Thurman, R.E., Stamatoyannopoulos, J.A., and Akey, J.M. (2012). Personal and population genomics of human regulatory variation. *Genome Res.* 22, 1689–1697.
- Voight, B.F., Kudaravalli, S., Wen, X., and Pritchard, J.K. (2006). A map of recent positive selection in the human genome. *PLoS Biol.* 4, e72.
- Wang, Y., Harvey, C.B., Pratt, W.S., Sams, V.R., Sarner, M., Rossi, M., Auricchio, S., and Swallow, D.M. (1995). The lactase persistence/non-persistence polymorphism is controlled by a cis-acting element. *Hum. Mol. Genet.* 4, 657–662.
- Wei, T., Gong, J., Rössle, S.C., Jamitzky, F., Heckl, W.M., and Stark, R.W. (2011). A leucine-rich repeat assembly approach for homology modeling of the human TLR5-10 and mouse TLR11-13 ectodomains. *J. Mol. Model.* 17, 27–36.

- Williamson, S.H., Hubisz, M.J., Clark, A.G., Payseur, B.A., Bustamante, C.D., and Nielsen, R. (2007). Localizing recent adaptive evolution in the human genome. *PLoS Genet.* 3, e90.
- Wong, S.H., Gochhait, S., Malhotra, D., Pettersson, F.H., Teo, Y.Y., Khor, C.C., Rautanen, A., Chapman, S.J., Mills, T.C., Srivastava, A., et al. (2010). Leprosy and the adaptation of human toll-like receptor 1. *PLoS Pathog.* 6, e1000979.
- Yi, X., Liang, Y., Huerta-Sanchez, E., Jin, X., Cuo, Z.X., Pool, J.E., Xu, X., Jiang, H., Vinckenbosch, N., Korneliussen, T.S., et al. (2010). Sequencing of 50 human exomes reveals adaptation to high altitude. *Science* 329, 75–78.
- Zhang, F.R., Huang, W., Chen, S.M., Sun, L.D., Liu, H., Li, Y., Cui, Y., Yan, X.X., Yang, H.T., Yang, R.D., et al. (2009). Genomewide association study of leprosy. *N. Engl. J. Med.* 361, 2609–2618.