

# Genome-wide detection and characterization of positive selection in human populations

Pardis C. Sabeti<sup>1\*</sup>, Patrick Varilly<sup>1\*</sup>, Ben Fry<sup>1</sup>, Jason Lohmueller<sup>1</sup>, Elizabeth Hostetter<sup>1</sup>, Chris Cotsapas<sup>1,2</sup>, Xiaohui Xie<sup>1</sup>, Elizabeth H. Byrne<sup>1</sup>, Steven A. McCarroll<sup>1,2</sup>, Rachele Gaudet<sup>3</sup>, Stephen F. Schaffner<sup>1</sup>, Eric S. Lander<sup>1,4,5,6</sup> & The International HapMap Consortium†

With the advent of dense maps of human genetic variation, it is now possible to detect positive natural selection across the human genome. Here we report an analysis of over 3 million polymorphisms from the International HapMap Project Phase 2 (HapMap2)<sup>1</sup>. We used 'long-range haplotype' methods, which were developed to identify alleles segregating in a population that have undergone recent selection<sup>2</sup>, and we also developed new methods that are based on cross-population comparisons to discover alleles that have swept to near-fixation within a population. The analysis reveals more than 300 strong candidate regions. Focusing on the strongest 22 regions, we develop a heuristic for scrutinizing these regions to identify candidate targets of selection. In a complementary analysis, we identify 26 non-synonymous, coding, single nucleotide polymorphisms showing regional evidence of positive selection. Examination of these candidates highlights three cases in which two genes in a common biological process have apparently undergone positive selection in the same population: *LARGE* and *DMD*, both related to infection by the Lassa virus<sup>3</sup>, in West Africa; *SLC24A5* and *SLC45A2*, both involved in skin pigmentation<sup>4,5</sup>, in Europe; and *EDAR* and *EDA2R*, both involved in development of hair follicles<sup>6</sup>, in Asia.

An increasing amount of information about genetic variation, together with new analytical methods, is making it possible to explore the recent evolutionary history of the human population. The first phase of the International Haplotype Map, including ~1 million single nucleotide polymorphisms (SNPs)<sup>7</sup>, allowed preliminary examination of natural selection in humans. Now, with the publication of the Phase 2 map (HapMap2)<sup>1</sup> in a companion paper, over 3 million SNPs have been genotyped in 420 chromosomes from three continents (120 European (CEU), 120 African (YRI) and 180 Asian from Japan and China (JPT + CHB)).

In our analysis of HapMap2, we first implemented two widely used tests that detect recent positive selection by finding common alleles carried on unusually long haplotypes<sup>2</sup>. The two, the Long-Range Haplotype (LRH)<sup>8</sup> and the integrated Haplotype Score (iHS)<sup>9</sup> tests, rely on the principle that, under positive selection, an allele may rise to high frequency rapidly enough that long-range association with nearby polymorphisms—the long-range haplotype<sup>8</sup>—will not have time to be eliminated by recombination. These tests control for local variation in recombination rates by comparing long haplotypes to other alleles at the same locus. As a result, they lose power as selected alleles approach fixation (100% frequency), because there are then

few alternative alleles in the population (Supplementary Fig. 2 and Supplementary Tables 1–2).

We next developed, evaluated and applied a new test, Cross Population Extended Haplotype Homozygosity (XP-EHH), to detect selective sweeps in which the selected allele has approached or achieved fixation in one population but remains polymorphic in the human population as a whole (Methods, and Supplementary Fig. 2 and Supplementary Tables 3–6). Related methods have recently also been described<sup>10–12</sup>.

Our analysis of recent positive selection, using the three methods, reveals more than 300 candidate regions (Supplementary Fig. 3 and Supplementary Table 7), 22 of which are above a threshold such that no similar events were found in 10 Gb of simulated neutrally evolving sequence (Methods). We focused on these 22 strongest signals (Table 1), which include two well-established cases, *SLC24A5* and *LCT*<sup>2,5,13</sup>, and 20 other regions with signals of similar strength.

The challenge is to sift through genetic variation in the candidate regions to identify the variants that were the targets of selection. Our candidate regions are large (mean length, 815 kb; maximum length, 3.5 Mb) and often contain multiple genes (median, 4; maximum, 15). A typical region harbours ~400–4,000 common SNPs (minor allele frequency >5%), of which roughly three-quarters are represented in current SNP databases and half were genotyped as part of HapMap2 (Supplementary Table 8).

We developed three criteria to help highlight potential targets of selection (Supplementary Fig. 1): (1) selected alleles detectable by our tests are likely to be derived (newly arisen), because long-haplotype tests have little power to detect selection on standing (pre-existing) variation<sup>14</sup>; we therefore focused on derived alleles, as identified by comparison to primate outgroups; (2) selected alleles are likely to be highly differentiated between populations, because recent selection is probably a local environmental adaptation<sup>2</sup>; we thus looked for alleles common in only the population(s) under selection; (3) selected alleles must have biological effects. On the basis of current knowledge, we therefore focused on non-synonymous coding SNPs and SNPs in evolutionarily conserved sequences. These criteria are intended as heuristics, not absolute requirements. Some targets of selection may not satisfy them, and some will not be in current SNP databases. Nonetheless, with ~50% of common SNPs in these populations genotyped in HapMap2, a search for causal variants is timely.

We applied the criteria to the regions containing *SLC24A5* and *LCT*, each of which already has a strong candidate gene, mutation and trait. At *SLC24A5*, the 600 kb region contains 914 genotyped

<sup>1</sup>Broad Institute of MIT and Harvard, Cambridge, Massachusetts 02139, USA. <sup>2</sup>Center for Human Genetic Research, Massachusetts General Hospital, Boston, Massachusetts 02114, USA. <sup>3</sup>Department of Molecular and Cellular Biology, Harvard University, Cambridge, Massachusetts 02138, USA. <sup>4</sup>Department of Biology, MIT, Cambridge, Massachusetts 02139, USA. <sup>5</sup>Whitehead Institute for Biomedical Research, Cambridge, Massachusetts 02142, USA. <sup>6</sup>Department of Systems Biology, Harvard Medical School, Boston, Massachusetts 02115, USA.

\*These authors contributed equally to this work.

†Lists of participants and affiliations appear at the end of the paper.

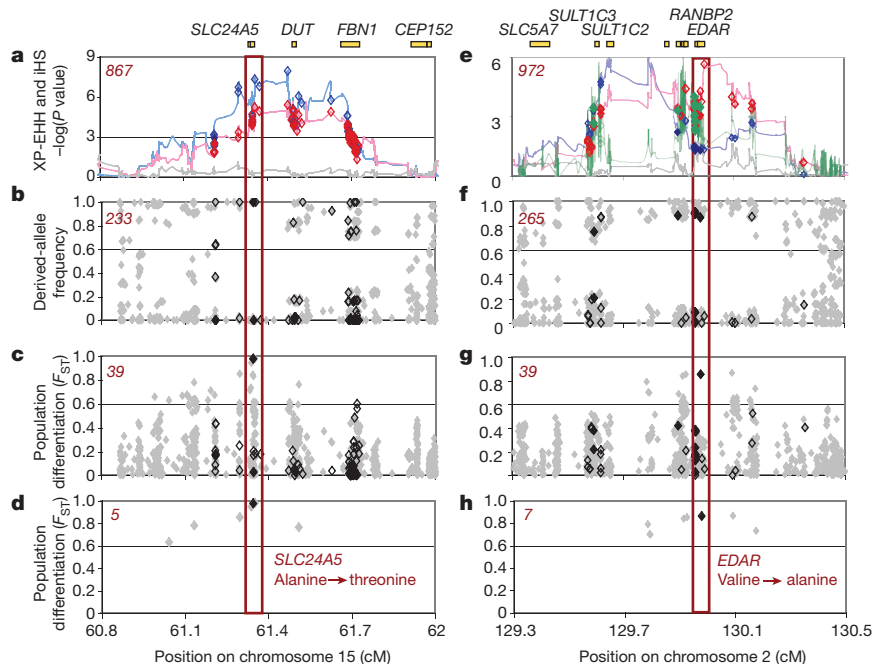
**Table 1 | The twenty-two strongest candidates for natural selection**

Region	Chr:position (MB, HG17)	Selected population	Long Haplotype Test	Size (Mb)	Total SNPs with Long Haplotype Signal	Subset of SNPs that fulfil criteria 1	Subset of SNPs that fulfil criteria 1 and 2	Subset of SNPs that fulfil criteria 1, 2 and 3	Genes at or near SNPs that fulfil all three criteria
1	chr1:166	CHB + JPT	LRH, iHS	0.4	92	39	30	2	<i>BLZF1, SLC19A2</i>
2	chr2:72.6	CHB + JPT	XP-EHH	0.8	732	250	0	0	
3	chr2:108.7	CHB + JPT	LRH, iHS, XP-EHH	1.0	972	265	7	1	<i>EDAR</i>
4	chr2:136.1	CEU	LRH, iHS, XP-EHH	2.4	1,213	282	24	3	<i>RAB3GAP1, R3HDM1, LCT</i>
5	chr2:177.9	CEU, CHB + JPT	LRH, iHS, XP-EHH	1.2	1,388	399	79	9	<i>PDE11A</i>
6	chr4:33.9	CEU, YRI, CHB + JPT	LRH, iHS	1.7	413	161	33	0	
7	chr4:42	CHB + JPT	LRH, iHS, XP-EHH	0.3	249	94	65	6	<i>SLC30A9</i>
8	chr4:159	CHB + JPT	LRH, iHS, XP-EHH	0.3	233	67	34	1	
9	chr10:3	CEU	LRH, iHS, XP-EHH	0.3	179	63	16	1	
10	chr10:22.7	CEU, CHB + JPT	XP-EHH	0.3	254	93	0	0	
11	chr10:55.7	CHB + JPT	LRH, iHS, XP-EHH	0.4	735	221	5	2	<i>PCDH15</i>
12	chr12:78.3	YRI	LRH, iHS	0.8	151	91	25	0	
13	chr15:46.4	CEU	XP-EHH	0.6	867	233	5	1	<i>SLC24A5</i>
14	chr15:61.8	CHB + JPT	XP-EHH	0.2	252	73	40	6	<i>HERC1</i>
15	chr16:64.3	CHB + JPT	XP-EHH	0.4	484	137	2	0	
16	chr16:74.3	CHB + JPT, YRI	LRH, iHS	0.6	55	35	28	3	<i>CHST5, ADAT1, KARS</i>
17	chr17:53.3	CHB + JPT	XP-EHH	0.2	143	41	0	0	
18	chr17:56.4	CEU	XP-EHH	0.4	290	98	26	3	<i>BCAS3</i>
19	chr19:43.5	YRI	LRH, iHS, XP-EHH	0.3	83	30	0	0	
20	chr22:32.5	YRI	LRH	0.4	318	188	35	3	<i>LARGE</i>
21	chr23:35.1	YRI	LRH, iHS	0.6	50	35	25	0	
22	chr23:63.5	YRI	LRH, iHS	3.5	13	3	1	0	
		Total SNPs		16.74	9,166	2,898	480	41	

Twenty-two regions were identified at a high threshold for significance (Methods), based on the LRH, iHS and/or XP-EHH test. Within these regions, we examined SNPs with the best evidence of being the target of selection on the basis of having a long haplotype signal, and by fulfilling three criteria: (1) being a high-frequency derived allele; (2) being differentiated between populations and common only in the selected population; and (3) being identified as functional by current annotation. Several candidate polymorphisms arise from the analysis including well-known *LCT* and *SLC24A5* (ref. 2), as well as intriguing new candidates.

SNPs. Applying filters progressively (Table 1 and Fig. 1a–d), we found that 867 SNPs are associated with the long-haplotype signal, of which 233 are high-frequency derived alleles, of which 12 are highly differentiated between populations, and of which only 5 are

common in Europe and rare in Asia and Africa. Among these five SNPs, there is only one implicated as functional by current knowledge; it has the strongest signal of positive selection and encodes the A111T polymorphism associated with pigment differences in

**Figure 1 | Localizing *SLC24A5* and *EDAR* signals of selection.**

**a–d, *SLC24A5*.** **a**, Strong evidence for positive selection in CEU samples at a chromosome 15 locus: XP-EHH between CEU and JPT + CHB (blue), CEU and YRI (red), and YRI and JPT + CHB (grey). SNPs are classified as having low probability (bordered diamonds) and high probability (filled diamonds) potential for function. SNPs were filtered to identify likely targets of selection on the basis of the frequency of derived alleles (**b**), differences between populations (**c**) and differences between populations for high-frequency derived alleles (less than 20% in non-selected populations) (**d**). The number of SNPs that passed each filter is given in the top left corner in red. The threonine to alanine candidate polymorphism in *SLC24A5* is the

clear outlier. **e–h, *EDAR*.** **e**, Similar evidence for positive selection in JPT + CHB at a chromosome 2 locus: XP-EHH between CEU and JPT + CHB (blue), between YRI and JPT + CHB (red), and between CEU and YRI (grey); iHS in JPT + CHB (green). A valine to alanine polymorphism in *EDAR* passes all filters: the frequency of derived alleles (**f**), differences between populations (**g**) and differences between populations for high-frequency derived alleles (less than 20% in non-selected populations) (**h**). Three other functional changes, a D→E change in *SULT1C2* and two SNPs associated with *RANBP2* expression (Methods), have also become common in the selected population.

humans and thought to be the target of positive selection<sup>5</sup>. Our criteria thus uniquely identify the expected allele.

At the *LCT* locus, we found similar degrees of filtration. Within the 2.4 Mb selective sweep, 24 polymorphisms fulfil the first two criteria (Table 1, and Supplementary Fig. 4), with the polymorphism thought to confer adult persistence of lactase among them. However, this SNP was only identified as functional after extensive study of the *LCT* gene<sup>15</sup>. Thus *LCT* shows both the utility and the limits of the heuristics.

Given the encouraging results for *SLC24A5* and *LCT*, we performed a similar analysis on all 22 candidate regions (Table 1). Filtering the 9,166 SNPs associated with the long-haplotype signal, we found that 480 satisfied the first two criteria. We identified 41 out of the 480 SNPs (0.2% of all SNPs genotyped in the regions) as possibly functional on the basis of a newly compiled database of polymorphisms in known coding elements, evolutionarily conserved elements and regulatory elements (Methods; B.F., unpublished), together containing ~ 5.5% of all known SNPs.

Eight of the forty-one SNPs encode non-synonymous changes (Table 1 and Supplementary Table 9). Apart from the well-known case of *SLC24A5*, they are found in *EDAR*, *PCDH15*, *ADAT1*, *KARS*, *HERC1*, *SLC30A9* and *BLZF1*. The remaining 33 potentially functional SNPs lie within conserved transcription factor motifs, introns, UTRs and other non-coding regions.

To identify additional candidates, we reversed the process by taking non-synonymous coding SNPs with highly differentiated high-frequency derived alleles; these SNPs comprise a tiny fraction of all SNPs and have a higher a priori probability of being targets of selection. Of the 15,816 non-synonymous SNPs in HapMap2, 281 (Supplementary Table 10) have both a high derived-allele frequency (frequency >50%) and clear differentiation between populations ( $F_{ST}$  is in the top 0.5 percentile). We examined these 281 SNPs to identify those embedded within long-range haplotypes<sup>16</sup>, and identified 26 putative cases of positive selection. These include the eight non-synonymous SNPs identified in the genome-wide analysis above.

Interestingly, analysis of the top regions and the non-synonymous SNPs together revealed three cases of two genes in the same pathway both having strong evidence of selection in a single population.

In the European sample, there is strong evidence for two genes already shown to be associated with skin pigment differences among humans. The first is *SLC24A5*, described above. We further examined the global distribution (Fig. 2) and the predicted effect on protein activity of the *SLC24A5* A111T polymorphism (Supplementary Fig. 5, 6). The second, *SLC45A2*, has an important role in pigmentation in zebrafish, mouse and horse<sup>4</sup>. An L374F substitution in *SLC45A2* is at 100% frequency in the European sample, but absent in the Asian and African samples. A recent association study has shown that the Phe-encoding allele is correlated with fair skin and non-black hair in Europeans<sup>4</sup>. Together, the data support *SLC45A2* as a target of positive selection in Europe<sup>10,17</sup>.

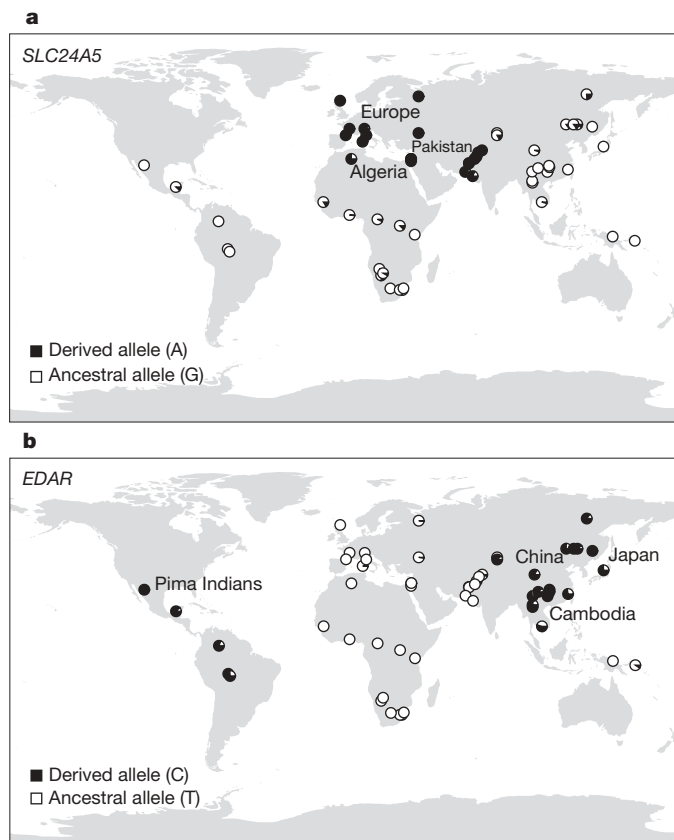
In the African sample (Yoruba in Ibadan, Nigeria), there is evidence of selection for two genes with well-documented biological links to the Lassa fever virus. The strongest signal in the genome, on the basis of the LRH test, resides within a 400 kb region that lies entirely within the gene *LARGE*. The *LARGE* protein is a glycosylase that post-translationally modifies  $\alpha$ -dystroglycan, the cellular receptor for Lassa fever virus (as well as other arenaviruses), and the modification has been shown to be critical for virus binding<sup>3</sup>. The virus name is derived from Lassa, Nigeria, where the disease is endemic, with 21% of the population showing signs of exposure<sup>18</sup>. We also noted that the *DMD* locus is on our larger candidate list of regions, with the signal of selection again in the Yoruba sample. *DMD* encodes a cytosolic adaptor protein that binds to  $\alpha$ -dystroglycan and is critical for its function. We hypothesize that Lassa fever created selective pressure at *LARGE* and *DMD*<sup>12</sup>. This hypothesis can be tested by correlating the geographical distribution of the selected haplotype

with endemicity of the Lassa virus, studying infection of genotyped cells *in vitro*, and searching for an association between the selected haplotype and clinical outcomes in infected patients.

In the Asian samples, we found evidence of selection for non-synonymous polymorphisms in two genes in the ectodysplasin (*EDA*) pathway, which is involved in development of hair, teeth and exocrine glands<sup>6</sup>. The genes are *EDAR* and *EDA2R*, which encode the key receptors for the ligands EDA A1 and EDA A2, respectively. Notably, the *EDA* signalling pathway has been shown to be under positive selection for loss of scales in multiple distinct populations of freshwater stickleback fish<sup>19</sup>. A mutation encoding a V370A substitution in *EDAR* is near fixation in Asia and absent in Europe and Africa (Fig. 1e–h). An R57K substitution in *EDA2R* has derived-allele frequencies of 100% in Asia, 70% in Europe and 0% in Africa.

The *EDAR* polymorphism is notable because it is highly differentiated between the Asian and other continental populations (the 3rd most differentiated among 15,816 non-synonymous SNPs), and also within Asian populations (in the top 1% of SNPs differentiated between the Japanese and Chinese HapMap samples). Genotyping of the *EDAR* polymorphism in the CEPH (Centre d'Etude du Polymorphisme Humain) global diversity panel<sup>20</sup> shows that it is at high but varying frequency throughout Asia and the Americas (for example, 100% in Pima Indians and in parts of China, and 73% in Japan) (Fig. 2, and Supplementary Fig. 7). Studying populations like the Japanese, in which the allele is still segregating, may provide clues to its biological significance.

*EDAR* has a central role in generation of the primary hair follicle pattern, and mutations in *EDAR* cause hypohidrotic ectodermal

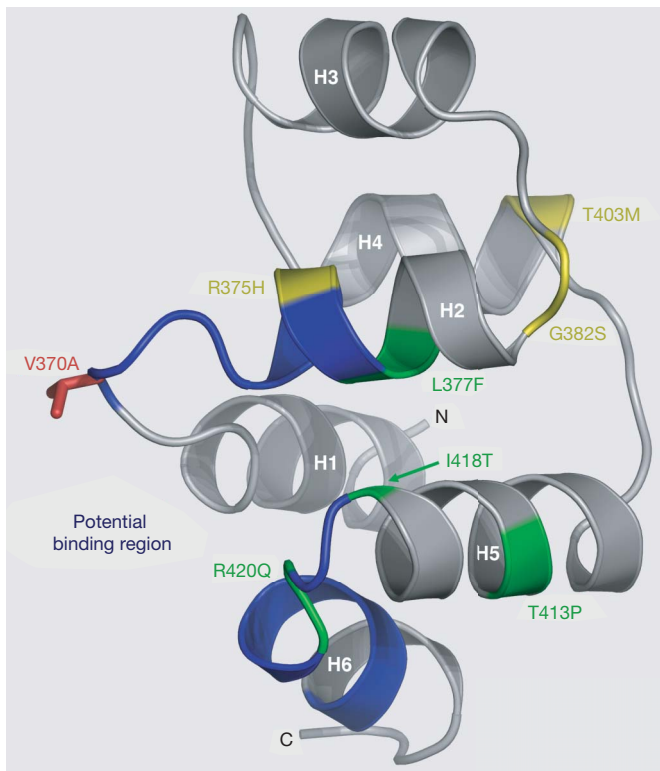


**Figure 2 | Global distribution of *SLC24A5* A111T and *EDAR* V370A.** Worldwide allele-frequency distributions for candidate polymorphisms with the strongest evidence for selection<sup>20</sup>. **a**, *SLC24A5* A111T is common in Europe, Northern Africa and Pakistan, but rare or absent elsewhere. **b**, *EDAR* V370A is common in Asia and the Americas, but absent in Europe and Africa.



dysplasia (HED) in humans and mice, characterized by defects in the development of hair, teeth and exocrine glands<sup>6</sup>. The V370A polymorphism, proposed to be the target of selection, lies within *EDAR*'s highly conserved death domain (Supplementary Fig. 8), the location of the majority of *EDAR* polymorphisms causing HED<sup>21</sup>. Our structural modelling predicts that the polymorphism lies within the binding site of the domain (Fig. 3).

Our analysis only scratches the surface of the recent selective history of the human genome. The results indicate that individual candidates may coalesce into pathways that reveal traits under selection, analogous to the alleles of multiple genes (for example, *HBB*, *G6PD* and *DARC*) that arose and spread in Africa and other tropical populations as a result of the partial protection they confer against malaria<sup>2,12</sup>. Such endeavours will be enhanced by continuing development of analytical methods to localize signals in candidate regions, generation of expanded data sets, advances in comparative genomics to define coding and regulatory regions, and biological follow-up of promising candidates. True understanding of the role of adaptive evolution will require collaboration across multiple disciplines, including molecular and structural biology, medical and population genetics, and history and anthropology.



**Figure 3 | Structural model of the EDAR death domain.** Ribbon representation of a homology model of the EDAR death domain (DD), based on the alignment of the EDAR DD amino acid sequence (EDAR residues 356–431), with multiple known DD structures. The helices are labelled H1 to H6. Residues in blue (the H1–H2 and H5–H6 loops, residues 370–376 and 419–425, respectively) correspond to the homologous residues in Tube that interact with Pelle in the Tube-DD–Pelle-DD structure<sup>24</sup>. These EDAR-DD residues therefore form a potential region of interaction with a DD-containing EDAR-interacting protein, such as EDARADD. The V370A polymorphic residue (red) is located prominently within this potential binding region in the H1–H2 loop. Seven of the thirteen known mis-sense mutations in EDAR that lead to hypohidrotic ectodermal dysplasia (HED) in humans are located in the EDAR-DD: the only four mutations in EDAR that lead to the dominant transmission of HED (green) and three recessive mutations (yellow)<sup>21</sup>. Four of these mutations, R375H, L377F, R420Q and I418T are located in the vicinity of the predicted interaction interface.

## METHODS SUMMARY

**Genotyping data.** Phase 2 of the International Haplotype Map (HapMap2) (www.hapmap.org) contains 3.1 million SNPs genotyped in 420 chromosomes in 3 continental populations (120 European (CEU), 120 African (YRI) and 180 Asian (JPT+CHB))<sup>1</sup>. We further genotyped our top HapMap2 functional candidates in the HGDR–CEPH Human Genome Diversity Cell Line Panel<sup>20</sup>.

**LRH, iHS and XP-EHH tests.** The Long-Range Haplotype (LRH), integrated Haplotype Score (iHS) and Cross Population EHH (XP-EHH) tests detect alleles that have risen to high frequency rapidly enough that long-range association with nearby polymorphisms—the long-range haplotype—has not been eroded by recombination; haplotype length is measured by the EHH<sup>8,9</sup>. The first two tests detect partial selective sweeps, whereas XP-EHH detects selected alleles that have risen to near fixation in one but not all populations. To evaluate the tests, we simulated genomic data for each HapMap population in a range of demographic scenarios—under neutral evolution and twenty scenarios of positive selection—developing the program Sweep (www.broad.mit.edu/mpg/sweep) for analysis. For our top candidates by the three tests, we tested for haplotype-specific recombination rates and copy-number polymorphisms, possible confounders.

**Localization.** We calculated  $F_{ST}$  and derived-allele frequency for all SNPs within the top candidate regions. We developed a database for those regions to annotate all potentially functional DNA changes (B.F., unpublished), including non-synonymous variants, variants disrupting predicted functional motifs, variants within regions of conservation in mammals and variants previously associated with human phenotypic differences, as well as synonymous, intronic and untranslated region variants.

**Structural model.** We generated a homology model of the EDAR death domain (DD) from available DD structures using Modeller 9v1 (ref. 22). The distribution of conserved residues, built using ConSurf<sup>23</sup> with an EDAR sequence alignment from 22 species, shows a bias to the protein core in helices H1, H2 and H5, supporting our model.

**Full Methods** and any associated references are available in the online version of the paper at www.nature.com/nature.

Received 8 August; accepted 13 September 2007.

1. The International HapMap Consortium. A second generation human haplotype map of over 3.1 million SNPs. *Nature* doi:10.1038/nature06258 (this issue).
2. Sabeti, P. C. *et al.* Positive natural selection in the human lineage. *Science* **312**, 1614–1620 (2006).
3. Kunz, S. *et al.* Posttranslational modification of  $\alpha$ -dystroglycan, the cellular receptor for arenaviruses, by the glycosyltransferase LARGE is critical for virus binding. *J. Virol.* **79**, 14282–14296 (2005).
4. Graf, J., Hodgson, R. & van Daal, A. Single nucleotide polymorphisms in the *MATP* gene are associated with normal human pigmentation variation. *Hum. Mutat.* **25**, 278–284 (2005).
5. Lamason, R. L. *et al.* SLC24A5, a putative cation exchanger, affects pigmentation in zebrafish and humans. *Science* **310**, 1782–1786 (2005).
6. Botchkarev, V. A. & Fessing, M. Y. Edar signaling in the control of hair follicle development. *J. Investig. Dermatol. Symp. Proc.* **10**, 247–251 (2005).
7. The International Haplotype Map Consortium. A haplotype map of the human genome. *Nature* **437**, 1299–1320 (2005).
8. Sabeti, P. C. *et al.* Detecting recent positive selection in the human genome from haplotype structure. *Nature* **419**, 832–837 (2002).
9. Voight, B. F., Kudaravalli, S., Wen, X. & Pritchard, J. K. A map of recent positive selection in the human genome. *PLoS Biol.* **4**, e72 (2006).
10. Kimura, R., Fujimoto, A., Tokunaga, K. & Ohashi, J. A practical genome scan for population-specific strong selective sweeps that have reached fixation. *PLoS ONE* **2**, e286 (2007).
11. Tang, K., Thornton, K. R. & Stoneking, M. A new approach for using genome scans to detect recent positive selection in the human genome. *PLoS Biol.* **5**, e171 (2007).
12. Williamson, S. H. *et al.* Localizing recent adaptive evolution in the human genome. *PLoS Genet.* **3**, e90 (2007).
13. Bersaglieri, T. *et al.* Genetic signatures of strong recent positive selection at the lactase gene. *Am. J. Hum. Genet.* **74**, 1111–1120 (2004).
14. Teshima, K. M., Coop, G. & Przeworski, M. How reliable are empirical genomic scans for selective sweeps? *PLoS Genet.* **16**, 702–712 (2006).
15. Kuokkanen, M. *et al.* Transcriptional regulation of the lactase–phlorizin hydrolase gene by polymorphisms associated with adult-type hypolactasia. *Gut* **52**, 647–652 (2003).
16. Miller, R. G. *Simultaneous statistical inference* XVI 299 (Springer, New York, 1981).
17. Soejima, M., Tachida, H., Ishida, T., Sano, A. & Koda, Y. Evidence for recent positive selection at the human *AIM1* locus in a European population. *Mol. Biol. Evol.* **23**, 179–188 (2006).
18. Richmond, J. K. & Baglolle, D. J. Lassa fever: epidemiology, clinical features, and social consequences. *Br. Med. J.* **327**, 1271–1275 (2003).
19. Colosimo, P. F. *et al.* Widespread parallel evolution in sticklebacks by repeated fixation of Ectodysplasin alleles. *Science* **307**, 1928–1933 (2005).

20. Rosenberg, N. A. *et al.* Genetic structure of human populations. *Science* **298**, 2381–2385 (2002).
21. Chassaing, N., Bourthoumieu, S., Cossee, M., Calvas, P. & Vincent, M. C. Mutations in *EDAR* account for one-quarter of non-*ED1*-related hypohidrotic ectodermal dysplasia. *Hum. Mutat.* **27**, 255–259 (2006).
22. Marti-Renom, M. A. *et al.* Comparative protein structure modeling of genes and genomes. *Annu. Rev. Biophys. Biomol. Struct.* **29**, 291–325 (2000).
23. Landau, M. *et al.* ConSurf 2005: the projection of evolutionary conservation scores of residues on protein structures. *Nucleic Acids Res.* **33**, W299–W302 (2005).
24. Xiao, T., Towb, P., Wasserman, S. A. & Sprang, S. R. Three-dimensional structure of a complex between the death domains of Pelle and Tube. *Cell* **99**, 545–555 (1999).

**Supplementary Information** is linked to the online version of the paper at [www.nature.com/nature](http://www.nature.com/nature).

**Acknowledgements** P.C.S. is funded by a Burroughs Wellcome Career Award in the Biomedical Sciences and has been funded by the Damon Runyon Cancer Fellowship and the L'Oréal for Women in Science Award. We thank A. Schier, B. Voight, R. Roberts, M. Kreiger, A. Abzhanov, D. Degusta, M. Burnette, E. Lieberman, M. Daly, D. Altschuler, D. Reich, D. Lieberman and I. Woods for helpful discussions on our analysis and results. We also thank L. Ziaugra, D. Tabbaa and T. Rachupka for experimental assistance. This work was funded in part by grants from the National Human Genome Research Institute (to E.S.L.) and from the Broad Institute of MIT and Harvard.

**Author Contributions** P.C.S., P.V., B.F. and E.S.L. initiated the project. P.V., B.F. and P.C.S. developed key software. P.C.S., P.V., B.F., S.F.S., J.L., E.H., C.C., X.X., E.B., S.A.McC. and R.G. performed analysis. P.C.S., E.B. and E.H. performed experiments. P.C.S., E.S.L., P.V. and S.F.S. wrote the manuscript.

**Author Information** Reprints and permissions information is available at [www.nature.com/reprints](http://www.nature.com/reprints). Correspondence and requests for materials should be addressed to P.C.S. ([pcdis@broad.mit.edu](mailto:pcdis@broad.mit.edu)).

**The International HapMap Consortium** (Participants are arranged by institution and then alphabetically within institutions except for Principal Investigators and Project Leaders, as indicated.)

**Genotyping centres:** **Perlegen Sciences** Kelly A. Frazer (Principal Investigator)<sup>1</sup>, Dennis G. Ballinger<sup>2</sup>, David R. Cox<sup>2</sup>, David A. Hinds<sup>2</sup>, Laura L. Stuve<sup>2</sup>; **Baylor College of Medicine and ParAllele BioScience** Richard A. Gibbs (Principal Investigator)<sup>3</sup>, John W. Belmont<sup>3</sup>, Andrew Boudreau<sup>4</sup>, Paul Hardenbol<sup>5</sup>, Suzanne M. Leal<sup>3</sup>, Shiran Pasternak<sup>6</sup>, David A. Wheeler<sup>3</sup>, Thomas D. Willis<sup>4</sup>, Fuli Yu<sup>7</sup>; **Beijing Genomics Institute** Huanming Yang (Principal Investigator)<sup>8</sup>, Changqing Zeng (Principal Investigator)<sup>8</sup>, Yang Gao<sup>8</sup>, Haoran Hu<sup>8</sup>, Weitao Hu<sup>8</sup>, Chaohua Li<sup>8</sup>, Wei Lin<sup>8</sup>, Siqi Liu<sup>8</sup>, Hao Pan<sup>8</sup>, Xiaoli Tang<sup>8</sup>, Jian Wang<sup>8</sup>, Wei Wang<sup>8</sup>, Jun Yu<sup>8</sup>, Bo Zhang<sup>8</sup>, Qingrun Zhang<sup>8</sup>, Hongbin Zhao<sup>8</sup>, Hui Zhao<sup>8</sup>, Jun Zhou<sup>8</sup>; **Broad Institute of Harvard and Massachusetts Institute of Technology** Stacey B. Gabriel (Project Leader)<sup>7</sup>, Rachel Barry<sup>7</sup>, Brendan Blumenstiel<sup>7</sup>, Amy Camargo<sup>7</sup>, Matthew Defelice<sup>7</sup>, Maura Faggart<sup>7</sup>, Mary Goyette<sup>7</sup>, Supriya Gupta<sup>7</sup>, Jamie Moore<sup>7</sup>, Huy Nguyen<sup>7</sup>, Robert C. Onofrio<sup>7</sup>, Melissa Parkin<sup>7</sup>, Jessica Roy<sup>7</sup>, Erich Stahl<sup>7</sup>, Ellen Winchester<sup>7</sup>, Liuda Ziaugra<sup>7</sup>, David Altschuler (Principal Investigator)<sup>7,9</sup>; **Chinese National Human Genome Center at Beijing** Yan Shen (Principal Investigator)<sup>10</sup>, Zhijian Yao<sup>10</sup>; **Chinese National Human Genome Center at Shanghai** Wei Huang (Principal Investigator)<sup>11</sup>, Xun Chu<sup>11</sup>, Yungang He<sup>11</sup>, Li Jin<sup>12</sup>, Yangfan Liu<sup>11</sup>, Yayun Shen<sup>11</sup>, Weiwei Sun<sup>11</sup>, Haifeng Wang<sup>11</sup>, Yi Wang<sup>11</sup>, Ying Wang<sup>11</sup>, Xiaoyan Xiong<sup>11</sup>, Liang Xu<sup>11</sup>; **Chinese University of Hong Kong** Mary M. Y. Waye (Principal Investigator)<sup>13</sup>, Stephen K. W. Tsui<sup>13</sup>; **Hong Kong University of Science and Technology** Hong Xue (Principal Investigator)<sup>14</sup>, J. Tze-Fei Wong<sup>14</sup>; **illumina** Luana M. Galver (Project Leader)<sup>15</sup>, Jian-Bing Fan<sup>15</sup>, Kevin Gunderson<sup>15</sup>, Sarah S. Murray<sup>1</sup>, Arnold R. Oliphant<sup>16</sup>, Mark S. Chee (Principal Investigator)<sup>17</sup>; **McGill University and Génomique Québec Innovation Centre** Alexandre Montpetit (Project Leader)<sup>18</sup>, Fanny Chagnon<sup>18</sup>, Vincent Ferretti<sup>18</sup>, Martin Leboeuf<sup>18</sup>, Jean-François Olivier<sup>4</sup>, Michael S. Phillips<sup>18</sup>, Stéphanie Roumy<sup>15</sup>, Clémentine Sallée<sup>19</sup>, Andrei Verner<sup>18</sup>, Thomas J. Hudson (Principal Investigator)<sup>20</sup>; **University of California at San Francisco and Washington University** Pui-Yan Kwok (Principal Investigator)<sup>21</sup>, Dongmei Cai<sup>21</sup>, Daniel C. Koboldt<sup>22</sup>, Raymond D. Miller<sup>22</sup>, Ludmila Pawlikowska<sup>21</sup>, Patricia Taillon-Miller<sup>22</sup>, Ming Xiao<sup>21</sup>; **University of Hong Kong** Lap-Chee Tsui (Principal Investigator)<sup>23</sup>, William Mak<sup>23</sup>, You Qiang Song<sup>23</sup>, Paul K. H. Tam<sup>23</sup>; **University of Tokyo and RIKEN** Yusuke Nakamura (Principal Investigator)<sup>24,25</sup>, Takahisa Kawaguchi<sup>25</sup>, Takuya Kitamoto<sup>25</sup>, Takashi Morizono<sup>25</sup>, Atsushi Nagashima<sup>25</sup>, Yoizo Ohnishi<sup>25</sup>, Akihiro Sekine<sup>25</sup>, Toshihiro Tanaka<sup>25</sup>, Tatsuhiko Tsunoda<sup>25</sup>; **Wellcome Trust Sanger Institute** Panos Deloukas (Project Leader)<sup>26</sup>, Christine P. Bird<sup>26</sup>, Marcos Delgado<sup>26</sup>, Emmanouil T. Dermizakis<sup>26</sup>, Rhian Gwilliam<sup>26</sup>, Sarah Hunt<sup>26</sup>, Jonathan Morrison<sup>27</sup>, Don Powell<sup>26</sup>, Barbara E. Stranger<sup>26</sup>, Pamela Whittaker<sup>26</sup>, David R. Bentley (Principal Investigator)<sup>28</sup>

**Analysis groups:** **Broad Institute** Mark J. Daly (Project Leader)<sup>7,9</sup>, Paul I. W. de Bakker<sup>7,9</sup>, Jeff Barrett<sup>7,9</sup>, Yves R. Chretien<sup>7</sup>, Julian Maller<sup>7,9</sup>, Steve McCarroll<sup>7,9</sup>, Nick Patterson<sup>7</sup>, Itzik Pe'er<sup>29</sup>, Alkes Price<sup>7</sup>, Shaun Purcell<sup>9</sup>, Daniel J. Richter<sup>7</sup>, Pardis Sabeti<sup>7</sup>, Richa Saxena<sup>7,9</sup>, Stephen F. Schaffner<sup>7</sup>, Pak C. Sham<sup>23</sup>, Patrick Varily<sup>7</sup>, David Altschuler

(Principal Investigator)<sup>7,9</sup>; **Cold Spring Harbor Laboratory** Lincoln D. Stein (Principal Investigator)<sup>6</sup>, Lalitha Krishnan<sup>6</sup>, Albert Vernon Smith<sup>6</sup>, Marcela K. Tello-Ruiz<sup>6</sup>, Gudmundur A. Thorisson<sup>30</sup>; **Johns Hopkins University School of Medicine** Aravinda Chakravarti (Principal Investigator)<sup>31</sup>, Peter E. Chen<sup>31</sup>, David J. Cutler<sup>31</sup>, Carl S. Kashuk<sup>31</sup>, Shin Lin<sup>31</sup>; **University of Michigan** Gonçalo R. Abecasis (Principal Investigator)<sup>32</sup>, Weihua Guan<sup>32</sup>, Yun Li<sup>32</sup>, Heather M. Munro<sup>33</sup>, Zhaohui Steve Qin<sup>32</sup>, Daryl J. Thomas<sup>34</sup>; **University of Oxford** Gilean McVean (Project Leader)<sup>35</sup>, Adam Auton<sup>35</sup>, Leonardo Bottolo<sup>35</sup>, Niall Cardin<sup>35</sup>, Susana Eyheramendy<sup>35</sup>, Colin Freeman<sup>35</sup>, Jonathan Marchini<sup>35</sup>, Simon Myers<sup>35</sup>, Chris Spencer<sup>7</sup>, Matthew Stephens<sup>36</sup>, Peter Donnelly (Principal Investigator)<sup>35</sup>; **University of Oxford, Wellcome Trust Centre for Human Genetics** Lon R. Cardon (Principal Investigator)<sup>37</sup>, Geraldine Clarke<sup>38</sup>, David M. Evans<sup>38</sup>, Andrew P. Morris<sup>38</sup>, Bruce S. Wei<sup>39</sup>; **RIKEN** Tatsuhiko Tsunoda (Principal Investigator)<sup>25</sup>, Todd A. Johnson<sup>25</sup>; **US National Institutes of Health** James C. Mullikin<sup>40</sup>; **US National Institutes of Health National Center for Biotechnology Information** Stephen T. Sherry<sup>41</sup>, Michael Feolo<sup>41</sup>, Andrew Skol<sup>42</sup>

**Community engagement/public consultation and sample collection groups:** **Beijing Normal University and Beijing Genomics Institute** Houcan Zhang<sup>43</sup>, Changqing Zeng<sup>8</sup>, Hui Zhao<sup>8</sup>; **Health Sciences University of Hokkaido, Eubios Ethics Institute, and Shinshu University** Ichiro Matsuda (Principal Investigator)<sup>44</sup>, Yoshimitsu Fukushima<sup>45</sup>, Darryl R. Mace<sup>46</sup>, Eiko Suda<sup>47</sup>; **Howard University and University of Ibadan** Charles N. Rotimi (Principal Investigator)<sup>48</sup>, Clement A. Adebamowo<sup>49</sup>, Ike Ajayi<sup>49</sup>, Toyin Anigwuo<sup>49</sup>, Patricia A. Marshall<sup>50</sup>, Chibuzor Nkwodimma<sup>49</sup>, Charmaine D. M. Royal<sup>48</sup>; **University of Utah** Mark F. Leppert (Principal Investigator)<sup>51</sup>, Missy Dixon<sup>51</sup>, Andy Peiffer<sup>51</sup>

**Ethical, legal and social issues:** **Chinese Academy of Social Sciences** Renzong Qiu<sup>52</sup>; **Genetic Interest Group** Alastair Kent<sup>53</sup>; **Kyoto University** Kazuto Kato<sup>54</sup>; **Nagasaki University** Norio Niikawa<sup>55</sup>; **University of Ibadan School of Medicine** Isaac F. Adewole<sup>49</sup>; **University of Montréal** Bartha M. Knoppers<sup>19</sup>; **University of Oklahoma** Morris W. Foster<sup>56</sup>; **Vanderbilt University** Ellen Wright Clayton<sup>57</sup>; **Wellcome Trust** Jessica Watkin<sup>58</sup>

**SNP discovery:** **Baylor College of Medicine** Richard A. Gibbs (Principal Investigator)<sup>3</sup>, John W. Belmont<sup>3</sup>, Donna Muzny<sup>3</sup>, Lynne Nazareth<sup>3</sup>, Erica Sodergren<sup>3</sup>, George M. Weinstock<sup>3</sup>, David A. Wheeler<sup>3</sup>, Imtaz Yakub<sup>3</sup>; **Broad Institute of Harvard and Massachusetts Institute of Technology** Stacey B. Gabriel (Project Leader)<sup>7</sup>, Robert C. Onofrio<sup>7</sup>, Daniel J. Richter<sup>7</sup>, Liuda Ziaugra<sup>7</sup>, Bruce W. Birren<sup>7</sup>, Mark J. Daly<sup>7,9</sup>, David Altschuler (Principal Investigator)<sup>7,9</sup>; **Washington University** Richard K. Wilson (Principal Investigator)<sup>59</sup>, Lucinda L. Fulton<sup>59</sup>; **Wellcome Trust Sanger Institute** Jane Rogers (Principal Investigator)<sup>26</sup>, John Burton<sup>26</sup>, Nigel P. Carter<sup>26</sup>, Christopher M. Clee<sup>26</sup>, Mark Griffiths<sup>26</sup>, Matthew C. Jones<sup>26</sup>, Kirsten McLay<sup>26</sup>, Robert W. Plumb<sup>26</sup>, Mark T. Ross<sup>26</sup>, Sarah K. Sims<sup>26</sup>, David L. Willey<sup>26</sup>

**Scientific management:** **Chinese Academy of Sciences** Zhu Chen<sup>60</sup>, Hua Han<sup>60</sup>, Le Kang<sup>60</sup>; **Genome Canada** Martin Godbout<sup>61</sup>, John C. Wallenburg<sup>62</sup>; **Génome Québec** Paul L'Archevêque<sup>63</sup>, Guy Bellemare<sup>63</sup>; **Japanese Ministry of Education, Culture, Sports, Science and Technology** Koji Saeki<sup>64</sup>; **Ministry of Science and Technology of the People's Republic of China** Hongguang Wang<sup>65</sup>, Daochang An<sup>65</sup>, Hongbo Fu<sup>65</sup>, Qing Li<sup>65</sup>, Zhen Wang<sup>65</sup>; **The Human Genetic Resource Administration of China** Renwu Wang<sup>66</sup>; **The SNP Consortium** Arthur L. Holden<sup>15</sup>; **US National Institutes of Health** Lisa D. Brooks<sup>67</sup>, Jean E. McEwen<sup>67</sup>, Mark S. Guyer<sup>67</sup>, Vivian Ota Wang<sup>67,68</sup>, Jane L. Peterson<sup>67</sup>, Michael Shi<sup>69</sup>, Jack Spiegel<sup>70</sup>, Lawrence M. Sung<sup>71</sup>, Lynn F. Zacharia<sup>67</sup>, Francis S. Collins<sup>72</sup>; **Wellcome Trust** Karen Kennedy<sup>61</sup>, Ruth Jamieson<sup>58</sup>, John Stewart<sup>58</sup>

<sup>1</sup>The Scripps Research Institute, 10550 North Torrey Pines Road MEM275, La Jolla, California 92037, USA. <sup>2</sup>Perlegen Sciences, 2021 Stierlin Court, Mountain View, California 94043, USA. <sup>3</sup>Baylor College of Medicine, Human Genome Sequencing Center, Department of Molecular and Human Genetics, 1 Baylor Plaza, Houston, Texas 77030, USA. <sup>4</sup>Affymetrix, 3420 Central Expressway, Santa Clara, California 95051, USA. <sup>5</sup>Pacific Biosciences, 1505 Adams Drive, Menlo Park, California 94025, USA. <sup>6</sup>Cold Spring Harbor Laboratory, 1 Bungtown Road, Cold Spring Harbor, New York 11724, USA. <sup>7</sup>The Broad Institute of Harvard and Massachusetts Institute of Technology, 1 Kendall Square, Cambridge, Massachusetts 02139, USA. <sup>8</sup>Beijing Genomics Institute, Chinese Academy of Sciences, Beijing 100300, China. <sup>9</sup>Massachusetts General Hospital and Harvard Medical School, Simches Research Center, 185 Cambridge Street, Boston, Massachusetts 02114, USA. <sup>10</sup>Chinese National Human Genome Center at Beijing, 3-707 N. Yongchang Road, Beijing Economic-Technological Development Area, Beijing 100176, China. <sup>11</sup>Chinese National Human Genome Center at Shanghai, 250 Bi Bo Road, Shanghai 201203, China. <sup>12</sup>Fudan University and CAS-MPG Partner Institute for Computational Biology, School of Life Sciences, SIBS, CAS, Shanghai, 201203, China. <sup>13</sup>The Chinese University of Hong Kong, Department of Biochemistry, The Croucher Laboratory for Human Genetics, 6/F Mong Man Wai Building, Shatin, Hong Kong. <sup>14</sup>Hong Kong University of Science and Technology, Department of Biochemistry and Applied Genomics Center, Clear Water Bay, Knowlton, Hong Kong. <sup>15</sup>illumina, 9885 Towne Centre Drive, San Diego, California 92121, USA. <sup>16</sup>Complete Genomics, 658 North Pastoria Avenue, Sunnyvale, California 94085, USA. <sup>17</sup>Prognosis Biosciences, 4215 Sorrento Valley Boulevard, Suite 105, San Diego, California 92121, USA. <sup>18</sup>McGill University and Génomique Québec Innovation Centre, 740 Dr Penfield Avenue, Montréal, Québec H3A 1A4, Canada. <sup>19</sup>University of Montréal, The Public Law Research Centre

(CRDP), PO Box 6128, Downtown Station, Montréal, Québec H3C 3J7, Canada.

<sup>20</sup>Ontario Institute for Cancer Research, MaRS Centre, South Tower, 101 College Street, Suite 500, Toronto, Ontario, M5G 1L7, Canada. <sup>21</sup>University of California, San Francisco, Cardiovascular Research Institute, 513 Parnassus Avenue, Box 0793, San Francisco, California 94143, USA. <sup>22</sup>Washington University School of Medicine, Department of Genetics, 660 S. Euclid Avenue, Box 8232, St Louis, Missouri 63110, USA. <sup>23</sup>University of Hong Kong, Genome Research Centre, 6/F, Laboratory Block, 21 Sassoon Road, Pokfulam, Hong Kong. <sup>24</sup>University of Tokyo, Institute of Medical Science, 4-6-1 Sirokanedai, Minatoku, Tokyo 108-8639, Japan. <sup>25</sup>RIKEN SNP Research Center, 1-7-22 Suehiro-cho, Tsurumi-ku Yokohama, Kanagawa 230-0045, Japan. <sup>26</sup>Wellcome Trust Sanger Institute, The Wellcome Trust Genome Campus, Hinxton, Cambridge CB10 1SA, UK. <sup>27</sup>University of Cambridge, Department of Oncology, Cambridge CB1 8RN, UK. <sup>28</sup>Solexa, Chesterford Research Park, Little Chesterford, nr Saffron Walden, Essex CB10 1XL, UK. <sup>29</sup>Columbia University, 500 West 120th Street, New York, New York 10027, USA. <sup>30</sup>University of Leicester, Department of Genetics, Leicester LE1 7RH, UK. <sup>31</sup>Johns Hopkins University School of Medicine, McKusick-Nathans Institute of Genetic Medicine, Broadway Research Building, Suite 579, 733 N. Broadway, Baltimore, Maryland 21205, USA. <sup>32</sup>University of Michigan, Center for Statistical Genetics, Department of Biostatistics, 1420 Washington Heights, Ann Arbor, Michigan 48109, USA. <sup>33</sup>International Epidemiology Institute, 1455 Research Boulevard, Suite 550, Rockville, Maryland 20850, USA. <sup>34</sup>Center for Biomolecular Science and Engineering, Engineering 2, Suite 501, Mail Stop CBSE/ITI, UC Santa Cruz, Santa Cruz, California 95064, USA. <sup>35</sup>University of Oxford, Department of Statistics, 1 South Parks Road, Oxford OX1 3TG, UK. <sup>36</sup>University of Chicago, Department of Statistics, 5734 S. University Avenue, Eckhart Hall, Room 126, Chicago, Illinois 60637, USA. <sup>37</sup>Fred Hutchinson Cancer Research Center, 1100 Fairview Avenue North, Seattle, Washington 98109, USA. <sup>38</sup>University of Oxford/Wellcome Trust Centre for Human Genetics, Roosevelt Drive, Oxford OX3 7BN, UK. <sup>39</sup>University of Washington Department of Biostatistics, Box 357232, Seattle, Washington 98195, USA. <sup>40</sup>US National Institutes of Health, National Human Genome Research Institute, 50 South Drive, Bethesda, Maryland 20892, USA. <sup>41</sup>US National Institutes of Health, National Library of Medicine, National Center for Biotechnology Information, 8600 Rockville Pike, Bethesda, Maryland 20894, USA. <sup>42</sup>University of Chicago, Department of Medicine, Section of Genetic Medicine, 5801 South Ellis, Chicago, Illinois 60637, USA. <sup>43</sup>Beijing Normal University, 19 Xijiekouwai Street, Beijing 100875, China. <sup>44</sup>Health Sciences University of Hokkaido, Ishikari Tobetsu Machi 1757, Hokkaido 061-0293, Japan. <sup>45</sup>Shinshu University School of Medicine, Department of Medical Genetics, Matsumoto 390-8621, Japan. <sup>46</sup>United

Nations Educational, Scientific and Cultural Organization (UNESCO Bangkok), 920 Sukhumwit Road, Prakanong, Bangkok 10110, Thailand. <sup>47</sup>University of Tsukuba, Eubios Ethics Institute, PO Box 125, Tsukuba Science City 305-8691, Japan. <sup>48</sup>Howard University, National Human Genome Center, 2216 6th Street, NW, Washington, District of Columbia 20059, USA. <sup>49</sup>University of Ibadan College of Medicine, Ibadan, Oyo State, Nigeria. <sup>50</sup>Case Western Reserve University School of Medicine, Department of Bioethics, 10900 Euclid Avenue, Cleveland, Ohio 44106, USA. <sup>51</sup>University of Utah, Eccles Institute of Human Genetics, Department of Human Genetics, 15 North 2030 East, Salt Lake City, Utah 84112, USA. <sup>52</sup>Chinese Academy of Social Sciences, Institute of Philosophy/Center for Applied Ethics, 2121, Building 9, Caoqiao Xinyuan 3 Qu, Beijing 100067, China. <sup>53</sup>Genetic Interest Group, 4D Leroy House, 436 Essex Road, London N130P, UK. <sup>54</sup>Kyoto University, Institute for Research in Humanities and Graduate School of Biostudies, Ushinomiya-cho, Sakyo-ku, Kyoto 606-8501, Japan. <sup>55</sup>Nagasaki University Graduate School of Biomedical Sciences, Department of Human Genetics, Sakamoto 1-12-4, Nagasaki 852-8523, Japan. <sup>56</sup>University of Oklahoma, Department of Anthropology, 455 W. Lindsey Street, Norman, Oklahoma 73019, USA. <sup>57</sup>Vanderbilt University, Center for Genetics and Health Policy, 507 Light Hall, Nashville, Tennessee 37232, USA. <sup>58</sup>Wellcome Trust, 215 Euston Road, London NW1 2BE, UK. <sup>59</sup>Washington University School of Medicine, Genome Sequencing Center, Box 8501, 4444 Forest Park Avenue, St Louis, Missouri 63108, USA. <sup>60</sup>Chinese Academy of Sciences, 52 Sanlihe Road, Beijing 100864, China. <sup>61</sup>Genome Canada, 150 Metcalfe Street, Suite 2100, Ottawa, Ontario K2P 1P1, Canada. <sup>62</sup>McGill University, Office of Technology Transfer, 3550 University Street, Montréal, Québec H3A 2A7, Canada. <sup>63</sup>Géome Québec, 630, boulevard René-Lévesque Ouest, Montréal, Québec H3B 1S6, Canada. <sup>64</sup>Ministry of Education, Culture, Sports, Science, and Technology, 3-2-2 Kasumigaseki, Chiyodaku, Tokyo 100-8959, Japan. <sup>65</sup>Ministry of Science and Technology of the People's Republic of China, 15 B. Fuxing Road, Beijing 100862, China. <sup>66</sup>The Human Genetic Resource Administration of China, b7, Zaojunmiao, Haidian District, Beijing 100081, China. <sup>67</sup>US National Institutes of Health, National Human Genome Research Institute, 5635 Fishers Lane, Bethesda, Maryland 20892, USA. <sup>68</sup>US National Institutes of Health, Office of Behavioral and Social Science Research, 31 Center Drive, Bethesda, Maryland 20892, USA. <sup>69</sup>Novartis Pharmaceuticals Corporation, Biomarker Development, One Health Plaza, East Hanover, New Jersey 07936, USA. <sup>70</sup>US National Institutes of Health, Office of Technology Transfer, 6011 Executive Boulevard, Rockville, Maryland 20852, USA. <sup>71</sup>University of Maryland School of Law, 500 W. Baltimore Street, Baltimore, Maryland 21201, USA. <sup>72</sup>US National Institutes of Health, National Human Genome Research Institute, 31 Center Drive, Bethesda, Maryland 20892, USA.



## METHODS

**Genotyping data.** The chromosomes examined in HapMap 2 were phased by the consortium using PHASE<sup>25</sup>.

The HGDR-CEPH Human Genome Diversity Cell Line Panel<sup>20</sup> consists of 1,051 individuals from 51 populations across the world. We obtained DNA for the panel from the Foundation Jean Dausset (CEPH) and genotyped our top functional candidates for selection in the panel.

**LRH, iHS, and XP-EHH tests.** The Long-Range Haplotype (LRH) and the integrated Haplotype Score (iHS) tests have been previously described<sup>8,9</sup> and our methods are given in Supplementary Methods.

EHH between two SNPs, A and B, is defined as the probability that two randomly chosen chromosomes are homozygous at all SNPs between A and B, inclusive<sup>6</sup>; it is usually calculated using a sample of chromosomes from a single population. Explicitly, if the  $N$  chromosomes in a sample form  $G$  homozygous groups, with each group  $i$  having  $n_i$  elements, EHH is defined as

$$\text{EHH} = \frac{\sum_{i=1}^G \binom{n_i}{2}}{\binom{N}{2}}$$

The XP-EHH test detects selective sweeps in which the selected allele has risen to high frequency or fixation in one population, but remains polymorphic in the human population as a whole; for this purpose it is more powerful than either iHS or LRH (Supplementary Fig. 2 and Supplementary Tables 3–6). XP-EHH uses cross-population comparison of haplotype lengths to control for local variation in recombination rates. Such cross-population comparison is complicated by the fact that haplotype lengths also depend on population history, such as bottlenecks and expansions<sup>26</sup>. The XP-EHH test normalizes for genome-wide differences in haplotype length between populations.

We define the XP-EHH test with respect to two populations, A and B, a given core SNP and a given direction (centromere distal or proximal). EHH is calculated for all SNPs in population A between the core SNP and X, and the value integrated with respect to genetic distance, with the result defined as  $I_A$ .  $I_B$  is defined analogously for population B. The statistic  $\ln(I_A/I_B)$  is then calculated; an unusually positive value suggests selection in population A, a negative value selection in B. For identifying outliers, the log-ratio is normalized to have zero mean and unit variance. Details are given in Supplementary Methods.

We developed a computer program, Sweep, to implement these tests (LRH, iHS and XP-EHH) for positive selection, (Supplementary Methods; www.broad.mit.edu/mpg/sweep). In identifying the 22 strongest candidate regions, we considered regions with signals in at least two of five tests (LRH, iHS and XP-EHH in the three pairwise comparisons among the three populations), as well as those that had the strongest signal for each individual test. With this threshold we found no events in 10 Gb of simulated neutrally evolving sequence. For the top candidates by the three tests, we have taken additional steps to rule out the effects of recombination rate variation and copy number polymorphisms (Supplementary Methods).

**Simulations and power calculations.** We simulated the evolution of 1 MB sections of 120 chromosomes from each of the three continental HapMap populations, using a previously validated demographic model<sup>27</sup>, under neutrality and under twenty scenarios of positive selection. We studied the effects of demography by further simulating recent bottlenecks with a range of intensity. Details of simulations and power calculations are given in Supplementary Methods.

**Functional annotation.** We developed an annotation database for our candidate regions to identify all DNA changes with potential functional consequence (B.F., unpublished). We first examined candidates most likely to be functional, including non-synonymous mutations, variants that disrupt predicted functional motifs (transcription factor motifs in conserved regions up to 10-kb 5' of known

genes and miRNA binding-site motifs in conserved 3' untranslated regions of known genes), and variations reported to be associated with human phenotypic differences. For the last category, we identified variations associated with a clinical state (for example, malaria resistance) by a review of the published literature and those associated with changes to gene expression in lymphoblastoid cell lines from the HapMap individuals. The annotation included insertion/deletion mutations of all sizes. We also examined candidates with lower probability of being functional, including synonymous, intronic and untranslated variations and those that occur within regions of conservation in mammalian species. These methods are described in greater detail in Supplementary Methods.

**Structural model of EDAR's death domain.** We generated a homology model for EDAR's death domain (DD) using six solved DD structures: p75 NGFR-DD, RAIDD-DD, Pelle-DD, FADD-DD, Fas-DD and IRAK4-DD<sup>24,28–32</sup>. We aligned the corresponding protein sequences using SALIGN<sup>33</sup>. We then added the amino acid sequence of EDAR's DD (residues 356–431) to this structural alignment using Modeller 9v1 (ref. 22). The resulting alignment was used as the input to Modeller 9v1 to build ten EDAR-DD structure models, and the best model was selected based on the Objective Function Score. Owing to the high DOPE scores in the H1–H2 loop we performed a loop refinement using Modeller9v1, significantly reducing the energy of this region. We further evaluated the model by examining the distribution of conserved residues using ConSurf<sup>3</sup> with an alignment of EDAR-DD sequences from 22 species. We observed a bias of conserved residues to the protein core in H1, H2 and H5, which supports our EDAR-DD model. To identify potential binding regions of EDAR-DD, we used LSQMAN<sup>34</sup> to superimpose the model to the Tube-DD–Pelle-DD complex structure<sup>24</sup>. The H1–H2 and H5–H6 loops of the EDAR-DD correspond to Tube residues interacting with Pelle, and H2–H3 and H4–H5 loops to Pelle residues interacting with Tube. We focused our analysis on the residues corresponding to the interacting region in Tube because our EDAR-DD model is most similar to Tube. Figures were generated with PyMOL<sup>35</sup>.

**Other analysis.** Description of methods for calculating  $F_{ST}$ , derived-allele frequency, alignment of the SLC24 amino acids, species alignments, conservation graphs, and estimation of the fraction of SNPs genotyped in HapMap2 and identified in dbSNP, are given in Supplementary Methods.

25. Stephens, M., Smith, N. J. & Donnelly, P. A new statistical method for haplotype reconstruction from population data. *Am. J. Hum. Genet.* **68**, 978–989 (2001).
26. Crawford, D. C. *et al.* Evidence for substantial fine-scale variation in recombination rates across the human genome. *Nature Genet.* **36**, 700–706 (2004).
27. Schaffner, S. F. *et al.* Calibrating a coalescent simulation of human genome sequence variation. *Genome Res.* **15**, 1576–1583 (2005).
28. Berglund, H. *et al.* The three-dimensional solution structure and dynamic properties of the human FADD death domain. *J. Mol. Biol.* **302**, 171–188 (2000).
29. Huang, B., Eberstadt, M., Olejniczak, E. T., Meadows, R. P. & Fesik, S. W. NMR structure and mutagenesis of the Fas (APO-1/CD95) death domain. *Nature* **384**, 638–641 (1996).
30. Lasker, M. V., Gajjar, M. M. & Nair, S. K. Cutting edge: molecular structure of the IL-1R-associated kinase-4 death domain and its implications for TLR signaling. *J. Immunol.* **175**, 4175–4179 (2005).
31. Liepinsh, E., Ilag, L. L., Otting, G. & Ibanez, C. F. NMR structure of the death domain of the p75 neurotrophin receptor. *EMBO J.* **16**, 4999–5005 (1997).
32. Park, H. H. & Wu, H. Crystal structure of RAIDD death domain implicates potential mechanism of PIDDosome assembly. *J. Mol. Biol.* **357**, 358–364 (2006).
33. Marti-Renom, M. A., Madhusudhan, M. S. & Sali, A. Alignment of protein sequences by their profiles. *Protein Sci.* **13**, 1071–1087 (2004).
34. Kleywegt, G. J. Use of non-crystallographic symmetry in protein structure refinement. *Acta Crystallogr. D* **52**, 842–857 (1996).
35. DeLano, W. L. *MacPyMOL: A PyMOL-based Molecular Graphics Application for MacOS X.* (DeLano Scientific LLC, Palo Alto, California, USA, 2007).