



HHS Public Access

Author manuscript

Nat Commun. Author manuscript; available in PMC 2015 September 27.

Published in final edited form as:

Nat Commun. ; 6: 6534. doi:10.1038/ncomms7534.

Human *gephyrin* is encompassed within giant functional noncoding yin-yang sequences

Sharlee Climer^{1,*}, Alan R. Templeton^{2,3,4}, and Weixiong Zhang^{1,3,5,*}

¹Department of Computer Science and Engineering, Washington University, St. Louis, MO 63130, USA

²Department of Biology, Washington University, St. Louis, MO 63130, USA

³Department of Genetics, Washington University, St. Louis, MO 63130, USA

⁴Department of Evolutionary and Environmental Biology, University of Haifa, Haifa 31905, Israel

⁵Institute for Systems Biology, Jiangnan University, Wuhan, Hubei 430056, China

Abstract

Gephyrin is a highly-conserved gene that is vital for the organization of proteins at inhibitory receptors, molybdenum cofactor biosynthesis, and other diverse functions. Its specific function is intricately regulated and its aberrant activities have been observed for a number of human diseases. Here we report a remarkable yin-yang haplotype pattern encompassing *gephyrin*. Yin-yang haplotypes arise when a stretch of DNA evolves to present two disparate forms that bear differing states for nucleotide variations along their lengths. The *gephyrin* yin-yang pair consists of 284 divergent nucleotide states and both variants vary drastically from their mutual ancestral haplotype, suggesting rapid evolution. Several independent lines of evidence indicate strong positive selection on the region and suggest these high-frequency haplotypes represent two distinct functional mechanisms. This discovery holds potential to deepen our understanding of variable human-specific regulation of *gephyrin* while providing clues for rapid evolutionary events and allelic migrations buried within human history.

INTRODUCTION

Gephyrin is a 93 kDa multi-functional protein that was named after the Greek word for 'bridge' due to its role in linking neurotransmitter receptors to the microtubule cytoskeleton. It binds polymerized tubulin with high affinity, likely due to a motif with high sequence similarities to the binding domains of MAP2 and tau^{1,2}. This protein dynamically provides a

Users may view, print, copy, and download text and data-mine the content in such documents, for the purposes of academic research, subject always to the full Conditions of use:http://www.nature.com/authors/editorial_policies/license.html#terms

*Corresponding author: climer@wustl.edu (S.C.); weixiong.zhang@wustl.edu (W.Z.).

AUTHOR CONTRIBUTIONS

S.C., A.R.T., and W.Z. conceived the project, designed the experiments, analyzed the data, interpreted the results, and wrote the manuscript. S.C. also performed the study.

COMPETING FINANCIAL INTERESTS

The authors declare no competing financial interests.

scaffold for clustering of proteins for both glycine and GABA-A receptors in inhibitory synapses, plays a crucial role in synapse formation and plasticity, and is believed to hold a central role in maintaining homeostatic excitation-inhibition balance³. Gephyrin has remarkably diverse functions. It associates with translation initiation machinery and has been implicated in the regulation of synaptic protein synthesis⁴. It also interacts with mammalian target of rapamycin (mTOR), a key protein for nutrient-sensitive cell cycle regulation, and has been shown to be required for downstream mTOR signaling⁵. Interestingly, gephyrin clustering at GABAergic synapses is increased by brain-derived neurotrophic factor (BDNF)-mediated mTOR activation and decreased by glycogen synthase kinase 3 beta (GSK3 β) phosphorylation⁶. Gephyrin is also indispensable for molybdenum cofactor (MoCo) biosynthesis as it is necessary for the insertion of molybdenum during this essential process³. MoCo deficiency leads to severe neurological damage and early childhood death. The fusion of an ancient function (MoCo biosynthesis) with an evolutionarily young function (neuroreceptor clustering) is believed to impact catalytic efficacy of MoCo synthesis by improving product-substrate channeling⁷. Finally, gephyrin was recently observed to localize within a ~600kDa cytoplasmic complex of unknown composition in non-neuronal cells, and it has been speculated that this complex might be involved in nutrient sensing, glucose metabolism, or aging, perhaps due to gephyrin's interactions with mTOR⁸.

Gephyrin's protein-coding regions are identical to the chimpanzee ortholog and are highly conserved across species. In contrast, regulation of this gene is highly variable. *Gephyrin* produces complex alternative splicing isoforms, which are crucial for its diverse functions, and at least eight of the 29 exons of this mosaic gene are subject to alternative splicing in species-, tissue-, cell-, and/or environmentally-specific manners^{1,9-13}. It is believed that the gephyrin scaffold in inhibitory synapses is a hexagonal lattice with twofold and three-fold symmetry and some alternative splicing isoforms disrupt this structure¹⁴. These alternate forms may provide a mechanism for plasticity and the dynamics of receptor anchoring by acting as dominant-negative variants which bind and remove receptors from synapses¹⁴. In concordance, MoCo biosynthesis activity is also isoform dependent, with various cassette insertions or deletions inactivating this synthesis¹⁵. For these reasons, unraveling the regulatory mechanisms is essential for elucidating and understanding gephyrin's dynamic and diverse activities and functions.

Markers within introns and in close genomic proximity are prominent candidates for regulatory elements, and the region encompassing *gephyrin* has been noted previously by two different groups. A 2.1 Mb region of homozygosity (ROH) in this location was discovered in 2010¹⁶. ROHs are correlated with linkage disequilibrium (LD) and have been observed to sometimes bear markedly disparate haplotypes¹⁷. In their 2010 paper, Curtis and Vine determined 20 genomic regions which had the largest number of subjects showing an ROH and studied the haplotypes of the nine single nucleotide polymorphisms (SNPs) at the center of each of these regions, observing that the haplotypes showed significant excess disparity, *i.e.* a tendency for pairs to simultaneously differ at multiple SNPs¹⁶. The term *yin-yang haplotypes* was coined to capture the polarity of such structures when a 24-SNP pattern for which two haplotypes with differing states at each site and a combined frequency of 0.50

was discovered by Zhang *et al.*¹⁸ Curtis and Vine noted that the ten most common haplotypes for the nine SNPs in the *gephyrin* region had a combined frequency of 0.67, indicating surprisingly little diversity of haplotypes. Interestingly, eight of these ten haplotypes yielded four pairs of yin-yang haplotypes, each of which bore different allelic states at all nine SNPs, indicating the haplotypes which did occur were remarkably different from each other.

In a 2012 study unrelated to yin-yang haplotypes, this region was identified by Park¹⁹ in a genome-wide scan of LD. This study identified an exceptionally strong LD block and discussed ‘extraordinary’ frequency spectra for all HapMap²⁰ populations in a 1 Mb region centered on intron 2 of *gephyrin*. Park concluded that the phenomenon could be due to a selective sweep and reviewed a number of selective pressure analyses, noting that this region had been included in supplementary materials by two of these studies^{21,22} and completely overlooked by the others. Park noted the uniqueness of this region, but the underlying yin-yang pattern went undetected.

In an exploration of genome-wide population data, we apply a recently developed method named BlocBuster²³ to SNP data for individuals in HapMap²⁰ populations and discover a high-frequency 284-SNP yin-yang haplotype pair embedded in noncoding regions within and surrounding *gephyrin*. Both haplotypes vary drastically from their mutual ancestral haplotype, yet they are highly conserved across global human populations, specifying two radically distinct evolutionary paths within a single genomic region. Furthermore, we report several independent lines of evidence indicating the identified yin and yang haplotypes are under selective pressure, thereby suggesting two distinct and functionally significant mechanisms underlie these regions.

RESULTS

***Gephyrin* is encompassed within a yin-yang haplotype pair**

We applied our BlocBuster method²³ (see Methods) to SNP data for unrelated individuals in four HapMap populations²⁰: CEU, CHB, JPT, and YRI. BlocBuster constructs networks that reveal haploid groups of SNP alleles that are inter-correlated, referred to as *blocs*. The results were highly consistent across all of the autosomal chromosomes, except chromosome 14, which had unusual network characteristics (Supplementary Note 1). We then applied BlocBuster to HapMap data for four different populations: GIH, LWK, MKK, and TSI²⁰, and again chromosome 14 was an outlier. A closer examination revealed the source of the anomalies - 73% of all of the edges in the first network were concentrated into a single bloc with 255 SNP alleles and 74% of the edges in the second network were concentrated into two blocs with 264 and 257 SNP alleles, respectively. The two blocs in the second network share 241 SNPs in common, with opposite alleles appearing in each bloc. Furthermore, these SNPs span the same genomic region as the bloc in the first network.

Overall, the three blocs found by the two analyses capture a single yin-yang haplotype pair. The three blocs possess 226 SNPs in common and span across 284 unique SNPs overall (Supplementary Dataset 1). We define this yin-yang pair using these 284 highly correlated SNPs. (See Supplementary Note 1 and Supplementary Fig. 1 for description of an additional

bloc corresponding to the yang haplotype for the first analysis.) This yin-yang pair is located on 14q23.3, encompassing *gephyrin* (*GPHN*) and extending beyond by ~300 kb upstream and downstream of *gephyrin* (Fig. 1). Interestingly, all of the divergent markers appear within introns, long non-coding RNA, or intergenic regions. As illustrated by the color coded bar above the first two columns of matrices in Fig. 1, few SNPs are downstream from *gephyrin* (2.7%, 3.4%, and 5.1% for the 255-, 264-, and 257-SNP blocs, respectively) and none lie within *MPP5*. About one-fifth of the SNPs lie upstream from *gephyrin* (19.2%, 18.6%, and 20.6%) and all three blocs include the same eight SNPs within the long non-coding RNA, *LINC00238*. Most of the SNPs lie within non-coding regions of *gephyrin* (78.0%, 78.0%, and 74.3%).

Due to the high proportion of heterozygotes within the Asian populations, we further interrogated these results using computationally phased haplotypes for the CHB and JPT populations provided by the HapMap Consortium. There are few yin-yang SNPs downstream from *gephyrin* and they are sparser and more variable than the other SNPs, so we omitted them from this analysis (see Methods). The available phased haplotypes did not include all of the yin-yang SNPs, and after removing SNPs with more than 5% missing data there remained 236 SNPs in phased haplotypes for 170 CHB+JPT individuals, for a total of 340 phased chromosomes. Fig. 2 shows the percentages of yin and yang SNP alleles found on each of the 340 phased chromosomes. These plots illustrate the prominence of the two divergent haplotypes and rarity of intermediate haplotypes.

Interleaving SNPs lying between the yin-yang SNPs generally have low minor allele frequencies (MAF), as shown in Fig. 3 and Supplementary Figs. 4–11. A close examination of the interleaving SNPs for the Asian populations indicate that a handful of individuals tend to possess most of the minor alleles (appearing in the high-resolution images as horizontal dotted lines across the yin-yang region of the matrix). Note that these individuals are not correlated with yin or yang haplotype status and consequently the variants are not likely to be hitchhiking with the yin or yang haplotypes.

These results indicate exceptionally high linkage amongst the 284 SNPs spanning more than 1 Mb and primarily located within noncoding regions of *gephyrin* and immediately upstream. Notably, two distinct haplotypes with differing states at all of the SNPs are unusually common and appear across global populations.

Conservation of yin-yang haplotypes within *Homo* populations

As shown in Fig. 1, the yin and yang haplotypes are prominent for all eleven HapMap populations, with combined frequencies ranging from 0.28 to 0.80. The pie charts in Fig. 1 indicate the frequencies of the yin and yang haplotypes and the white regions represent the portion of partial haplotypes with one or more alleles that do not conform to an entire yin or yang pattern. The percentages of homozygotes and heterozygotes are listed on the right of each matrix. The two European-ancestry populations, CEU and TSI, have large percentages of yin homozygotes. Three African populations, LWK, MKK, and YRI, have high frequencies of the yang haplotypes and possess a recombination block near the end of the haplotypes, while individuals with African ancestry in Southwest USA (ASW) have a yang frequency of 0.07 and a shorter recombination block.

The East and South Asian populations (CHB, CHD, GIH and JPT) exhibit the strongest mix of yin and yang haplotypes. Every one of these four populations exhibit frequencies of at least 0.25 for each of the yin and yang haplotypes and the combined frequencies for the CHB and JPT populations reach 0.76 and 0.80, respectively. The Chinese in Metropolitan Denver, Colorado, USA (CHD) are similar to the Han Chinese in Beijing, China (CHB), although there is some recombination near the start of the haplotypes for the CHD. The GIH, with ancestry from the Indian subcontinent, possess frequencies that are similar to the East Asian populations, albeit with decreased yang homozygotes and increased haplotype diversity.

The 1000 Genomes Project²⁴ includes genotype data for 2,504 individuals from 26 global populations, representing each major human ancestry. Although imputed data are included in these files, we built a BlocBuster network to test the robustness of the results found for the HapMap data, as described in Supplementary Note 2. The yin and yang haplotypes are pronounced for these individuals (Supplementary Fig. 2), thereby supporting the HapMap results.

Ancestral alleles for the 284 yin-yang SNPs were determined by comparing human and chimpanzee DNA (see Methods) and are shown in Fig. 4. Both the yin and yang haplotypes are significantly different from the ancestral haplotype, sharing only 51.4% and 48.6% identity by state (IBS), respectively. The macaque, orangutan, and chimpanzee haplotypes are also shown in Fig. 4 and are generally similar to the ancestral haplotype.

The available Neandertal and Denisovan data also predominantly match the ancestral alleles. Fig. 4 displays 15 SNP alleles for three Neandertal and the single individual available from the Denisovan fossil site^{25,26} (Neand/Denis) that have been typed on the Affymetrix HuOrigin array²⁷. The SNP ascertainment approach for the HuOrigin array had a bias for SNPs with matching Denisovan and chimpanzee alleles (see Methods). As shown in Fig. 4, all but one of the 15 matches the chimpanzee and ancestral alleles. In all, 11 of the 15 SNP alleles, including the derived allele, match the yin haplotype. Also shown in Fig. 4 are high-coverage genotypes for the Denisovan individual²⁸. In contrast to the Neand/Denis data, all 125 SNPs match the yang haplotype. Although there is no apparent reason to expect a bias in these data, 95.2% of the alleles are IBS with both the ancestral and chimpanzee alleles. This is unexpected as less than half of the yang alleles are IBS with the ancestral alleles.

Overall, while the yin-yang genotypic patterns are not conserved across species outside the *Homo* genus, they are highly conserved across the HapMap populations, with combined frequencies ranging from 0.28 to 0.80 for the pair, as detailed in Fig. 1.

Selection for the yin and yang haplotypes

Several lines of evidence suggest the yin and yang haplotypes are under strong positive selection and bear functional importance. First, a series of diverse statistical tests for selection indicate positive selection for the region, as shown in Fig. 5. The left panel of the figure shows the results for four selection tests computed over four HapMap populations. The right panel shows results for selection tests computed over the 1000 Genomes Project data (released April 2012). The topmost plot on the right represents a selective sweep scan

on Neandertal vs. human polymorphisms, followed by rank scores for 13 tests for selection²⁹ (see Methods). Both panels include results from Fay and Wu's H test³⁰. This test was specifically designed to distinguish between positive selection and background selection by utilizing data from outgroup species. As shown in the figure, the yin-yang interval has a statistically significant H value. Taken together, these results indicate strong positive selection within the yin-yang region.

It is worth noting that Nielsen *et al.* found that *gephyrin* showed no evidence for positive selection (p-value = 1.0)³¹ in the coding regions of the gene. Their calculations were specifically based upon the ratio of nonsynonymous to synonymous mutations within coding regions. In view of the strong selection pressure in the host genomic region, but not upon *gephyrin* exons, it follows that selection pressures may be acting upon functional elements within noncoding regions.

Second, the size, composition, and geographic distribution of yin and yang haplotypes indicate rapid evolution suggestive of strong positive selection. While the 284 identified SNPs have zero IBS between the yin and yang pair, the appearance of these haplotypes across eleven diverse populations must be identity by descent (IBD) for each haplotype, due to the identical states of hundreds of SNP alleles. Recall that the yin haplotype is prominent in European populations, yang is prominent in African populations, and Asian populations have nearly equal proportions. This observation, along with the assumption that the haplotypes are IBD, suggests that the Asian occurrences arose via gene flow or admixture. It follows that more than one hundred nucleotide mutations became fixed for *each* of the two haplotypes after their split from each other and prior to their migration to Asia. Such rapid evolution is indicative of strong selection. Surprisingly, these mutations remain generally fixed in these haplotypes in modern populations and all of the intermediate haplotypes that arose between the initial split and fixed states have low frequencies or have disappeared entirely.

Third, the unusual recombination patterns in this region support selection favoring the yin and yang haplotypes. A close examination of Figs. 1, 2, and 3 suggests that recombinants comprised of both a yin and yang parental haplotype are generally rare, particularly within *gephyrin* and upstream from this gene. Such a recombinant would appear as a horizontal bar comprised of blocks that are two different colors in Fig. 1. As shown in Fig. 2, nine of the 340 CHB and JPT haplotypes have between 10% and 90% yin/yang compositions; six of these represent yin-yang recombinants and three represent intermediate yin or yang haplotypes with more than 10% mutational variations. Indeed, the prevalence of each of the distinct yin and yang haplotypes, despite strong coexistence and recombination opportunities, indicates very low recombination events between yin and yang haplotypes. However, as shown in Fig. 6, previous analyses of this region provided by the HapMap Consortium (<http://hapmap.ncbi.nlm.nih.gov/downloads/recombination/>) reported moderate recombination within the region, including an estimated recombination rate of 9.2 cM/Mb at rs10133120 in the 5' end of *gephyrin*. Taken together, these results strongly suggest that recombinants comprised of two yin haplotypes and/or recombinants comprised of two yang haplotypes are more prevalent than recombinants merging yin and yang haplotypes together.

This observation suggests that yin and yang haplotypes may have been favorably selected over merged yin and yang recombinants.

DISCUSSION

It has been estimated that 5% of the human genome is under selection, yet only about 1% of the genome is protein-coding³², indicating that selection acts upon more noncoding than coding regions. Furthermore, transcription is pervasive and about 70% to 90% of the human genome is transcribed, producing a vast array of noncoding RNA³³. Some long noncoding RNA (lncRNA) have been documented to play critical regulatory roles. For example, the X-inactivate-specific transcript (XIST) is vital for inactivating the X chromosome for females by directly binding an epigenetic complex. Closer to protein-coding regions, untranslated regions (UTRs) contain the internal ribosome entry sites and riboswitches that participate in regulation of expression as well as alternative splicing³⁴. Furthermore, 3'UTRs host binding sites for microRNAs that inhibit translation³⁵. Intronic regions can provide noncoding RNA and are also involved in alternative splicing and transcription regulation³⁶. Alternatively, transcription of antisense strands can produce noncoding RNAs involved in a variety of biological roles³⁷.

The protein-coding regions of gephyrin are highly conserved and its diverse roles are accomplished via regulatory variations. Noncoding elements within its introns and upstream are prime candidates for such regulatory control. Importantly, aberrant regulation of this gene has been associated with a host of complex diseases. Dysfunction in the regulation of gephyrin expression levels and/or isoform production has been implicated for Alzheimer's disease (AD)³⁸⁻⁴⁰, epilepsy^{9,41,42}, autism¹⁰, schizophrenia^{10,43}, hyperekplexia¹³, and chorein deficiency⁴⁴. Gephyrin levels are significantly reduced in AD brains³⁹, and the normally strong correlations between gephyrin production and the abundances of the six most common GABA subunits is corrupted in AD brains³⁸. It has also been observed that abnormal accumulations of low-molecular weight gephyrin plaques overlap beta-amyloid plaques⁴⁰. Epilepsy is characterized by abnormal excessive excitatory neuronal activities, and dysfunction of inhibitory neurons and/or down-regulation of inhibitory circuits may be the underlying cause⁴¹. Gephyrin plays a vital role in inhibitory circuits. Both reduced levels of gephyrin production, as well as the appearance of aberrant gephyrin isoforms, have been observed in epileptogenesis^{9,41,42}. In individuals lacking gephyrin mutations, four aberrant gephyrin isoforms with missing exons have been observed to arise due to cellular stress. These isoforms display dominant negative effects on normal gephyrin in epileptogenesis⁹. Other alternative isoforms have been identified as risk factors for autism and schizophrenia and may also act as dominant-negative variants¹⁰. Athanasiu et al. conducted a GWAS of schizophrenia in Norwegian and European samples and tabulated 32 SNPs in the human genome with the most significant associations⁴³. Seven of the 32 are amongst the yin-yang SNPs, specifically: rs1952070, rs6573695, rs17247749, rs17836572, rs1885198, rs6573706, and rs7154017. Overall, the associations of gephyrin regulation with a half-dozen complex diseases strongly motivate the need to understand the genetic machinery driving the diverse manifestations of this highly conserved gene.

We present a remarkably long yin-yang haplotype pair spanning the noncoding regions of *gephyrin*. This genetic phenomenon is more than an order of magnitude larger than any previously reported yin-yang pair and is prevalent across eleven global human populations. Despite the conservation of these haplotypes across human populations, both are highly dissimilar to their common ancestral haplotype, suggesting they are the result of two divergent human-specific evolutionary paths. We advance this hypothesis by reporting several independent lines of evidence supporting selection for the two haplotypes. Taken together, this research lays the groundwork for a deep understanding of the regulatory control of *gephyrin*.

It is not clear how this genetic anomaly arose. Mutation and recombination have created vast amounts of haplotype diversity in many species, including humans. Previous reports have suggested that human-specific traits evolved primarily due to positive selection in noncoding regions involved in the regulation of genes⁴⁵⁻⁴⁷. The most eminent of these characteristics is the human brain, with its increased size and enhanced cognition, and it has been demonstrated that selection acting upon noncoding regions is predominantly associated with neural development, whereas selection acting upon protein coding regions is associated with immunity, olfaction, and male reproduction⁴⁷. In short, it is viable to expect that human-specific adaptations of *gephyrin* are due to evolution of regulatory mechanisms lying within noncoding regions, particularly those in close proximity.

A key question follows: why would two extremely divergent paths arise during such adaptation? One possibility is a chromosomal inversion resulting with a lack of recombination between the original and inverted variant. In such an event, the original and inverted haplotypes would evolve independently. Strong positive selection could drive the evolution of a single high-frequency haplotype for each group. Several systematic searches for inversions have been conducted over the human genome⁴⁸⁻⁵⁰. The most recent investigation mapped 6.1 million clones to distinct genomic positions for eight HapMap individuals (4 YRI, 2 CEU, 1 CHB, and 1 JPT) and identified 224 inversions⁵⁰. One of these is a 31.1 kb inversion in *FUT8*, which is 763 kb upstream from *gephyrin*. However, none of the three studies identified an inversion in the yin-yang region.

Another possible impetus for this pattern could be incompatible mutations. Suppose two independent mutations each possess a selective advantage individually, but the combination of the two mutations reduces fitness. For example, each of the mutations could increase the expression of a particular gene in a beneficial manner, but together they produce deleteriously high expression. Selection would favor haplotypes possessing either mutation, and recombinants possessing both or neither mutation would become rare. Over time the two haplotypes bearing each of the original mutations would evolve in distinct manners.

At least one other alternate mechanism could have led to the extreme divergence of the yin and yang haplotypes: convergent evolution in isolated ancient populations followed by gene flow⁵¹. Opportunities for such events have been common throughout human history. For example, recent sequencing of fossil DNA has led to an estimate that modern non-African populations may possess approximately 1.5% to 2.1% Neandertal DNA⁵². DNA related to the single individual found at Denisova is also found in modern island Southeast Asia and

Oceania populations, with modern Papuans possessing 6% of their DNA closely related to the Denisovan individual's DNA²⁸. As shown in Fig. 4, the Neandertal and Denisovan genotypes are highly similar to the ancestral haplotype. However, in addition to the small number of Neandertal genotypes, another weakness of this analysis is that the currently available data is based on few individuals. Increased sample size, increased marker density, and further investigations, such as comparisons with nuclear DNA from the 300,000-year-old hominins from Sima de los Huesos⁵³ when it becomes available, are needed to determine the likelihood that ancient admixture lies at the root of this yin-yang.

All of the described hypothetical mechanisms are likely to exhibit *differential* recombination, as is observed for the *gephyrin* yin-yang pair. The recombination rate amongst yin haplotypes and the rate amongst yang haplotypes appear substantially higher than the rate between yin and yang parental haplotypes. Selection is likely to be the strongest for chromosomal inversions as a recombination event between yin and yang haplotypes results with too few or too many copies of genes upstream and downstream from the crossover point and general abolition of a gene spanning this point. The existence of recombinants, including those with crossover points within *gephyrin*, casts doubt that an inversion underlies this anomaly. On a different note, a test for differential recombination might prove to be a valuable tool for assessing functionality of other yin-yang haplotype pairs previously identified and those to be mapped in the coming years. In general, if the yin and yang haplotypes are not functional, this type of differential recombination across coexisting haplotypes would be improbable.

The forces that produced this phenomenon, as well as the biological implications of its presence, invite exploration of an evolutionary 'road less traveled' that produced two highly divergent, and uniquely human, genetic patterns intricately interwoven with the conserved protein-coding regions of *gephyrin*. These results solicit new questions and provide material for hypotheses generation. Several avenues of future research have appealing potential, a couple of which are highlighted below.

With regard to *gephyrin* in particular, deep sequencing of the yin-yang region for ancient and modern populations could be valuable for discerning molecular-level function as well as providing insights into the historical journeys of the haplotypes. Also, testing for associations between yin-yang status and various phenotypes could provide valuable knowledge. Candidate phenotypes include transcript isoforms, variations of gene expression, and susceptibilities to complex diseases such as epilepsy, autism, and schizophrenia, which have previously shown to be associated with distinct isoforms of *gephyrin*^{9,10}. It should be noted that the use of animal models in previous studies of *gephyrin* might have been confounded and misleading as both the yin and yang haplotypes are uniquely human.

More generally, mapping of additional yin-yang haplotypes within the human genome, and other genomes of interest, may pinpoint genetic mechanisms underlying convergent pathways and/or expose regions undergoing rapid evolution. Additionally, when combined with geographic distributions, these patterns may provide distinguishable flags for understanding the histories of individuals and populations. Importantly, they may capture valuable features of an individual's genetic background and their susceptibility to complex

traits, perhaps aiding personalized medicine. Looking forward, in addition to increasing our understanding of the human-specific regulation of a vitally important gene, this haplotype pair may serve as a model for studying yin-yang haplotypes and their biological implications for human health and development.

METHODS

HapMap data

HapMap bulk data were downloaded from <http://hapmap.ncbi.nlm.nih.gov/>. Release HapMap r28, nr.b36 dated 18-Aug-2010 files were downloaded from directory/downloads/genotypes/2010-08_phaseII+III/forward/. Some of the individuals were related, as tabulated here: http://hapmap.ncbi.nlm.nih.gov/downloads/samples_individuals_relationships_w_pops_121708.txt. Data for the children were removed from the datasets, leaving presumably unrelated individuals. For each analysis, the SNPs that were common for all four populations were determined. Then these data were cleaned to reduce the quantity of missing genotypes as follows. First, the SNPs with at least 50% missing data were removed, then the individuals with at least 50% missing data were removed, and finally SNPs with at least 10% missing data were removed. The remaining individuals also had no more than 10% missing data.

In the first analysis, data for four populations were considered: Northern and western European ancestry (CEU), Han Chinese in Beijing (CHB), Japanese in Tokyo (JPT), and Yoruba in Ibadan, Nigeria (YRI). After removing the children and cleaning, the final data consisted of 1,115,561 autosomal SNPs for 112 CEU, 137 CHB, 113 JPT, and 116 YRI; a total of 478 individuals.

In the second analysis, data for four different populations were used: Gujarati Indians living in Houston, USA, with at least three grandparents from Gujarat (the northwest region of the Indian subcontinent) (GIH); Luhya in Webuye, Kenya (LWK); Maasai in Kinyawa, Kenya (MKK); and Toscani in Italia (TSI). After removing children and cleaning the data, the final data consisted of 1,242,039 autosomal SNPs for 101 GIH, 110 LWK, 143 MKK, and 102 TSI; a total of 456 individuals.

The three remaining HapMap populations were used to further validate the yin-yang haplotype pair: African ancestry in Southwest USA (ASW); Chinese in Metropolitan Denver, Colorado (CHD); and Mexican ancestry in Los Angeles, California (MEX). For each population, the genotypes for the 284 SNPs were extracted, when available, and haplotype frequencies were computed. The genotypes were also plotted for visual inspection (Fig. 1). All of the processed datasets can be obtained by contacting the first author.

1000 Genomes data

Chromosome 14 data were downloaded from the 1000 Genomes Project website at <ftp://ftp-trace.ncbi.nlm.nih.gov/1000genomes/ftp/release/20130502/> on Nov. 17, 2014. File ALL.chr14.phase3_shapeit2_mvncall_integrated_v5.20130502.genotypes.vcf.gz, with the last modification noted on Sept. 17, 2014, was obtained. A total of 2,504 individuals were genotyped. All markers between 66974125 and 67648525 (GRCh37 coordinates) were

extracted, yielding 13,992 markers in the yin-yang region. We extracted the 13,564 biallelic SNPs within this set.

Neandertal and Denisova data

One Neandertal and two Denisovan datasets were utilized. The Neand/Denis data for three Vindija Neandertal²⁵ and one individual from the Denisovan fossil site²⁶ were downloaded from ftp://ftp.cephb.fr/hgdp_supp10/Harvard_HGDP-CEPH/annotation.txt. The Affymetrix HuOrigin array²⁷ was used and included 15 SNPs from the yin-yang haplotypes. The data file includes the numbers of high quality reads for each allele. Only one of the 15 SNPs had more than one nucleotide state detected for all of the Neandertal and Denisovan reads. SNP rs6573754 (AX-50160621) had one 'A' and five 'G's for the Denisovan individual, and three 'G's for the Neandertal. The 'G' allele is shown for Neand/Denis in Fig. 4 of the main paper and Supplementary Dataset 1.

Panel 13 of the SNP ascertainment for the HuOrigin array only included SNPs for which the Denisovan allele matched the chimpanzee allele, as this policy facilitated validations²⁷. This panel accounted for 20.2% of the original 750,184 SNPs selected, presenting some bias when comparing the 15 Neand/Denis alleles with chimpanzee and ancestral alleles.

The second set of Denisovan data was downloaded from UCSC's Table Browser website (<http://genome.ucsc.edu/cgi-bin/hgTables>) by selecting the 'Denisova Assembly and Analysis' group and 'Denisova Variants' track from the Human GRCh37/hg19 assembly. The genetic material was drawn from the inner portion of the phalanx of the same individual represented in the Neand/Denis data. A single-stranded library preparation method was used to produce the high-coverage sequence²⁸. These data included 125 of the yin-yang haplotype SNPs, three of which were amongst the 15 SNPs in the Neand/Denis data.

BlocBuster

BlocBuster is a network approach that utilizes a multi-faceted, allele-oriented correlation measure^{23,54}. Briefly, we developed the approach with an aim to identify combinations of correlated alleles that are subjected to genetic heterogeneity. The correlation metric, CCC, is customized for genotype data and appreciates heterogeneity by evaluating four distinct correlations that retain independence between different types of pair-wise correlations. This specification of correlation types is retained in an allele-specific network construction which increases the network infrastructure yet maintains high efficiency. We determined the CCC threshold using the default method of setting the number of edges in the network equal to the number of SNPs. After preprocessing and cleaning the data, there were 36,542 SNPs in the CEU, CHB, JPT, YRI chromosome 14 dataset and 40,820 SNPs in the GIH, LWK, MKK, TSI chromosome 14 dataset and each of the networks contained the corresponding number of edges, representing the most significant CCC correlations for each analysis. Consequently, the average degree of each node in each of the networks was one. The significance of this correlation threshold was tested using permutation trials²³. After the networks were constructed, groups of nodes that were connected by edges were readily identified as they were completely isolated from each other. Each of these groups of connected nodes, referred to as blocs, represent a haploid pattern of inter-correlated SNP

alleles. The entire pattern of SNP alleles for each of these blocs was tested for possession by each individual. Our open-source code is available at www.blocbuster.org or by contacting the first author.

Determination of ancestral allelic similarities

Ancestral alleles were compiled from NCBI's dbSNP webpage (<http://www.ncbi.nlm.nih.gov/projects/SNP/>). These alleles were supplied by Dr. Jim Mullikin of the National Human Genome Research Institute (NHGRI) and were determined by comparing human and chimpanzee DNA⁵⁵. A complete list of the alleles for the 284 unique SNPs is supplied in Supplementary Dataset 1. Haplotype similarities were measured by tallying the numbers of markers that were identical by state (IBS), a simple yet accurate metric⁵⁶.

Selection tests

The selection test results were drawn from three sources. First, the Haplotter²² website (<http://haplotter.uchicago.edu/>) was utilized to plot results for four statistics: iHs, H, D, and F_{ST} , for four HapMap populations (CEU, CHB, JPT, and YRI) over a 5 Mb region centered on *gephyrin*. Voight *et al.*'s integrated haplotype score (iHs) is based upon an integration of the extended haplotype homozygosity (EHH) statistic²² and is designed to capture very recent positive selection. Fay and Wu's H statistic detects the effects of hitchhiking on the frequency spectrum as a function of recombination rate³⁰. Tajima's D statistic tests the neutral mutation hypothesis based on the relationship between the average number of nucleotide differences and the number of segregating sites⁵⁷. The fixation index, F_{ST} , is based on Wright's measure of population differentiation.

Second, the 1000 Genomes Selection Browser 1.0²⁹ was used to plot the results for a number of statistical tests computed over the 1000 Genomes Project data (<http://www.1000genomes.org/>) for CEU, CHB, and YRI populations. These resequencing data yield higher density information than the original HapMap data and remove most of the SNP ascertainment bias, making them valuable for summary statistics. We included the rank scores, which were computed using an outlier approach based on sorted genome-wide scores²⁹. Peaks in the plots represent regions under positive selection. Some of the methods were modified by Pybus *et al.*²⁹ and are marked in the following with an asterisk. Three families of statistical tests were included: Allele Frequency Spectrum, Linkage Disequilibrium Structure, and Population Differentiation. The Allele Frequency Spectrum family included Tajima's D (Taj_D)⁵⁷, Fay and Wu's H (FayWu_H)³⁰, Fu and Li's D (FuLi_D)⁵⁸, Fu and Li's F (FuLi_F)⁵⁸, and Ramos-Onsins and Rozas' R2 (R2)⁵⁹. The Linkage Disequilibrium Structure family included Sabeti *et al.*'s XP-EHH* (XPEHH)⁶⁰, Sabeti *et al.*'s EHH_average* (EHH)⁶¹, Nei's Dh (Dh)⁶², and Kelly's ZnS (ZnS)⁶³. The Population Differentiation family included Weir and Cockerham's pairwise F_{st} (Fst)⁶⁴, Chen *et al.*'s XP-CLR (XPCLR)⁶⁵, Hofer *et al.*'s absolute DAF (absDAF)⁶⁶, and Hofer *et al.*'s standard DAF (DAF)⁶⁶.

Third, we used selection statistics generated by Nielsen *et al.* which were determined by comparing synonymous and nonsynonymous mutations within coding regions³¹. More

specifically, the ratio of nonsynonymous substitutions per nonsynonymous site to synonymous substitutions per synonymous site was tested against the neutral null hypothesis of the ratio being one.

Determination of protein conservation

The conservation of the gephyrin protein across species was determined using UCSD Signaling Gateway (<http://www.signaling-gateway.org/molecule/>).

Haplotype data

The haplotypes for the combined CHB and JPT individuals were identified as follows. First, the SNPs that lie within gephyrin or upstream from gephyrin were extracted from the full HapMap data (positions 65,893,425 to 66,709,924 from HapMap r28, nr.b36). After removing the five individuals (2 CHB and 3 JPT) with excessive missing data, the SNPs with more than 5% missing data were discarded, leaving 326 SNPs. Next, the phased haplotypes from the same region for the JPT+CHB individuals were downloaded from the HapMap website (http://hapmap.ncbi.nlm.nih.gov/downloads/phasing/2009-02_phaseIII/HapMap3_r2/). These haplotypes had been inferred using PHASE^{67,68} and included 170 individuals from the combined CHB and JPT data. We discarded all SNPs that had been identified as having more than 5% missing values in the original genotype data, leaving a total of 303 phased sites.

Haplotype composition plots

The haplotype composition plots were constructed using the phased haplotypes for the CHB +JPT populations. These haplotypes were computationally inferred using PHASE⁶⁷. Of the 303 phased sites, 236 represented divergent yin-yang SNPs and the haplotypes comprised of these 236 SNP alleles were extracted for the 170 individuals. For each of the 340 phased chromosomes, the percentages of SNP alleles matching the yin and yang haplotypes, respectively, were computed.

Genotype heat maps

The genotype values for SNPs in the yin-yang region were plotted for visual inspection (Figs. 1 and 3). Individuals (rows) were reordered in order to place similar individuals near each other. We utilized our rearrangement clustering method, TSP+k⁶⁹ for this reordering. Briefly, the genotype values for the SNPs for each pattern were extracted from the data and converted to an instance of the Traveling Salesman Problem (TSP)⁷⁰ in which each individual was represented as a city. We inserted a dummy city to provide a natural break to the circular TSP tour and determined the ordering of the cities using an iterated Lin-Kernighan local search as implemented by Applegate, Bixby, Chvatal, and Cook in the Concorde package (<http://www.math.uwaterloo.ca/tsp/concorde/index.html>). The individuals were reordered using this solution and the genotypes were color-encoded with dark blue, light blue, red, and white representing homozygote for the identified allele, heterozygote, homozygote for the alternate allele, and missing data, respectively.

Supplementary Material

Refer to Web version on PubMed Central for supplementary material.

Acknowledgments

We thank Carlos Cruchaga, Michael Garvin, Alison Goate, Christina Gurnett, Cynthia C. Vigueira, and Patrick Vigueira for helpful discussions and David Reich for supplying Neandertal and Denisova data. HapMap bulk data were downloaded from <http://hapmap.ncbi.nlm.nih.gov/>. This work was supported by the National Institutes of Health [grant numbers P50-GM65509, RC1-AR058681, R01-GM086412, R01-GM100364], the National Science Foundation [grant number DBI-0743797], and the municipal government of Wuhan, Hubei, China (grant number 2014070504020241 and the Talent Development Program).

References

1. Ramming M, et al. Diversity and phylogeny of gephyrin: tissue-specific splice variants, gene structure, and sequence similarities to molybdenum cofactor-synthesizing and cytoskeleton-associated proteins. *Proceedings of the National Academy of Sciences of the United States of America*. 2000; 97:10266–71. [PubMed: 10963686]
2. Kirsch J, et al. The 93-kDa glycine receptor-associated protein binds to tubulin. *J Biol Chem*. 1991; 266:22242–22245. [PubMed: 1657993]
3. Tyagarajan SK, Fritschy JM. Gephyrin: a master regulator of neuronal function? *Nature reviews. Neuroscience*. 2014; 15:141–56. [PubMed: 24552784]
4. Sertie AL, de Alencastro G, De Paula VJ, Passos-Bueno MR. Collybistin and gephyrin are novel components of the eukaryotic translation initiation factor 3 complex. *BMC research notes*. 2010; 3:242. [PubMed: 20858277]
5. Sabatini DM, et al. Interaction of RAFT1 with gephyrin required for rapamycin-sensitive signaling. *Science (New York, NY)*. 1999; 284:1161–4.
6. Wuchter J, et al. A comprehensive small interfering RNA screen identifies signaling pathways required for gephyrin clustering. *The Journal of neuroscience_ : the official journal of the Society for Neuroscience*. 2012; 32:14821–34. [PubMed: 23077067]
7. Belaidi AA, Schwarz G. Metal insertion into the molybdenum cofactor: product-substrate channelling demonstrates the functional origin of domain fusion in gephyrin. *The Biochemical journal*. 2013; 450:149–57. [PubMed: 23163752]
8. Nawrotzki R, Islinger M, Vogel I, Völkl A, Kirsch J. Expression and subcellular distribution of gephyrin in non-neuronal tissues and cells. *Histochemistry and cell biology*. 2012; 137:471–82. [PubMed: 22270318]
9. Förster B, et al. Irregular RNA splicing curtails postsynaptic gephyrin in the cornu ammonis of patients with epilepsy. *Brain_ : a journal of neurology*. 2010; 133:3778–94. [PubMed: 21071388]
10. Lionel AC, et al. Rare exonic deletions implicate the synaptic organizer Gephyrin (GPHN) in risk for autism, schizophrenia and seizures. *Human molecular genetics*. 2013; 22:2055–66. [PubMed: 23393157]
11. Herweg J, Schwarz G. Splice-specific glycine receptor binding, folding, and phosphorylation of the scaffolding protein gephyrin. *The Journal of biological chemistry*. 2012; 287:12645–56. [PubMed: 22351777]
12. Meier J, Grantyn R. A gephyrin-related mechanism restraining glycine receptor anchoring at GABAergic synapses. *The Journal of neuroscience_ : the official journal of the Society for Neuroscience*. 2004; 24:1398–405. [PubMed: 14960612]
13. Rees MI, et al. Isoform heterogeneity of the human gephyrin gene (GPHN), binding domains to the glycine receptor, and mutation analysis in hyperekplexia. *The Journal of biological chemistry*. 2003; 278:24688–96. [PubMed: 12684523]
14. Bedet C, et al. Regulation of gephyrin assembly and glycine receptor synaptic stability. *The Journal of biological chemistry*. 2006; 281:30046–56. [PubMed: 16882665]

15. Smolinsky B, Eichler SA, Buchmeier S, Meier JC, Schwarz G. Splice-specific functions of gephyrin in molybdenum cofactor biosynthesis. *The Journal of biological chemistry*. 2008; 283:17370–9. [PubMed: 18411266]
16. Curtis D, Vine AE. Yin yang haplotypes revisited - long, disparate haplotypes observed in European populations in regions of increased homozygosity. *Human heredity*. 2010; 69:184–92. [PubMed: 20203523]
17. Curtis D, Vine AE, Knight J. Study of regions of extended homozygosity provides a powerful method to explore haplotype structure of human populations. *Annals of human genetics*. 2008; 72:261–78. [PubMed: 18205893]
18. Zhang J, Rowe WL, Clark AG, Buetow KH. Genomewide distribution of high-frequency, completely mismatching SNP haplotype pairs observed to be common across human populations. *American journal of human genetics*. 2003; 73:1073–81. [PubMed: 14560401]
19. Park L. Linkage disequilibrium decay and past population history in the human genome. *PloS one*. 2012; 7:e46603. [PubMed: 23056365]
20. Altshuler DM, et al. Integrating common and rare genetic variation in diverse human populations. *Nature*. 2010; 467:52–8. [PubMed: 20811451]
21. Williamson SH, et al. Localizing recent adaptive evolution in the human genome. *PLoS genetics*. 2007; 3:e90. [PubMed: 17542651]
22. Voight BF, Kudaravalli S, Wen X, Pritchard JK. A map of recent positive selection in the human genome. *PLoS biology*. 2006; 4:e72. [PubMed: 16494531]
23. Climer S, Templeton AR, Zhang W. Allele-Specific Network Reveals Combinatorial Interaction That Transcends Small Effects in Psoriasis GWAS. *PLoS Computational Biology*. 2014; 10:e1003766. [PubMed: 25233071]
24. Abecasis GR, et al. A map of human genome variation from population-scale sequencing. *Nature*. 2010; 467:1061–73. [PubMed: 20981092]
25. Green RE, et al. A draft sequence of the Neandertal genome. *Science (New York, NY)*. 2010; 328:710–22.
26. Reich D, et al. Genetic history of an archaic hominin group from Denisova Cave in Siberia. *Nature*. 2010; 468:1053–60. [PubMed: 21179161]
27. Patterson N, et al. Ancient admixture in human history. *Genetics*. 2012; 192:1065–93. [PubMed: 22960212]
28. Meyer M, et al. A high-coverage genome sequence from an archaic Denisovan individual. *Science (New York, NY)*. 2012; 338:222–6.
29. Pybus M, et al. 1000 Genomes Selection Browser 1.0: a genome browser dedicated to signatures of natural selection in modern humans. *Nucleic acids research*. 2014; 42:D903–9. [PubMed: 24275494]
30. Fay JC, Wu CI. Hitchhiking under positive Darwinian selection. *Genetics*. 2000; 155:1405–13. [PubMed: 10880498]
31. Nielsen R, et al. A scan for positively selected genes in the genomes of humans and chimpanzees. *PLoS biology*. 2005; 3:e170. [PubMed: 15869325]
32. Waterston RH, et al. Initial sequencing and comparative analysis of the mouse genome. *Nature*. 2002; 420:520–62. [PubMed: 12466850]
33. Lee JT. Epigenetic regulation by long noncoding RNAs. *Science (New York, NY)*. 2012; 338:1435–9.
34. Ray PS, et al. A stress-responsive RNA switch regulates VEGFA expression. *Nature*. 2009; 457:915–9. [PubMed: 19098893]
35. Bartel DP. MicroRNAs: target recognition and regulatory functions. *Cell*. 2009; 136:215–33. [PubMed: 19167326]
36. Barrett LW, Fletcher S, Wilton SD. Regulation of eukaryotic gene expression by the untranslated gene regions and other non-coding elements. *Cellular and molecular life sciences_*: CMLS. 2012; 69:3613–34. [PubMed: 22538991]
37. Faghihi MA, Wahlestedt C. Regulatory roles of natural antisense transcripts. *Nature reviews. Molecular cell biology*. 2009; 10:637–43. [PubMed: 19638999]

38. Limon A, Reyes-Ruiz JM, Mileli R. Loss of functional GABA(A) receptors in the Alzheimer diseased brain. *Proceedings of the National Academy of Sciences of the United States of America*. 2012; 109:10071–6. [PubMed: 22691495]
39. Agarwal S, Tannenberg RK, Dodd PR. Reduced Expression of the Inhibitory Synapse Scaffolding Protein Gephyrin in Alzheimer’s Disease. *Journal of Alzheimer’s disease*. 2008; 14:313–321. [PubMed: 18599957]
40. Hales CM, et al. Abnormal gephyrin immunoreactivity associated with Alzheimer disease pathologic changes. *Journal of neuropathology and experimental neurology*. 2013; 72:1009–15. [PubMed: 24128675]
41. González MI. The possible role of GABAA receptors and gephyrin in epileptogenesis. *Frontiers in cellular neuroscience*. 2013; 7:113. [PubMed: 23885234]
42. Fang M, et al. Downregulation of gephyrin in temporal lobe epilepsy neurons in humans and a rat model. *Synapse (New York, NY)*. 2011; 65:1006–14.
43. Athanasiu L, et al. Gene variants associated with schizophrenia in a Norwegian genome-wide study are replicated in a large European cohort. *Journal of psychiatric research*. 2010; 44:748–53. [PubMed: 20185149]
44. Kurano Y, et al. Chorein deficiency leads to upregulation of gephyrin and GABA(A) receptor. *Biochemical and biophysical research communications*. 2006; 351:438–42. [PubMed: 17070500]
45. Prabhakar S, Noonan JP, Pääbo S, Rubin EM. Accelerated evolution of conserved noncoding sequences in humans. *Science (New York, NY)*. 2006; 314:786.
46. O’Bleness M, Searles VB, Varki A, Gagneux P, Sikela JM. Evolution of genetic and genomic features unique to the human lineage. *Nature reviews Genetics*. 2012; 13:853–66.
47. Haygood R, Babbitt CC, Fedrigo O, Wray GA. Contrasts between adaptive coding and noncoding changes during human evolution. *Proceedings of the National Academy of Sciences of the United States of America*. 2010; 107:7853–7. [PubMed: 20385805]
48. Korb J, et al. Paired-end mapping reveals extensive structural variation in the human genome. *Science (New York, NY)*. 2007; 318:420–6.
49. Tuzun E, et al. Fine-scale structural variation of the human genome. *Nature genetics*. 2005; 37:727–32. [PubMed: 15895083]
50. Kidd JM, et al. Mapping and sequencing of structural variation from eight human genomes. *Nature*. 2008; 453:56–64. [PubMed: 18451855]
51. Templeton A. Out of Africa again and again. *Nature*. 2002; 416:45–51. [PubMed: 11882887]
52. Prüfer K, et al. The complete genome sequence of a Neanderthal from the Altai Mountains. *Nature*. 2014; 505:43–9. [PubMed: 24352235]
53. Meyer M, et al. A mitochondrial genome sequence of a hominin from Sima de los Huesos. *Nature*. 2013; 505:403–6. [PubMed: 24305051]
54. Climer S, Yang W, de Las Fuentes L, Dávila-Román VG, Gu CC. A Custom Correlation Coefficient (CCC) Approach for Fast Identification of Multi-SNP Association Patterns in Genome-Wide SNPs Data. *Genetic epidemiology*. 2014; 38:610–621. [PubMed: 25168954]
55. Spencer CCA, et al. The influence of recombination on human genetic diversity. *PLoS genetics*. 2006; 2:e148. [PubMed: 17044736]
56. Marquard V, Beckmann L, Bermejo JL, Fischer C, Chang-Claude J. Comparison of measures for haplotype similarity. *BMC proceedings*. 2007; 1 (Suppl 1):S128. [PubMed: 18466470]
57. Tajima F. Statistical method for testing the neutral mutation hypothesis by DNA polymorphism. *Genetics*. 1989; 123:585–95. [PubMed: 2513255]
58. Fu YX, Li WH. Statistical tests of neutrality of mutations. *Genetics*. 1993; 133:693–709. [PubMed: 8454210]
59. Ramos-Onsins SE, Rozas J. Statistical Properties of New Neutrality Tests Against Population Growth. *Molecular Biology and Evolution*. 2002; 19:2092–2100. [PubMed: 12446801]
60. Sabeti PC, et al. Genome-wide detection and characterization of positive selection in human populations. *Nature*. 2007; 449:913–8. [PubMed: 17943131]
61. Sabeti PC, et al. Detecting recent positive selection in the human genome from haplotype structure. *Nature*. 2002; 419:832–7. [PubMed: 12397357]

62. Marks J. Molecular evolutionary genetics. *American Journal of Physical Anthropology*. 1988; 75:428–429.
63. Kelly JK. A test of neutrality based on interlocus associations. *Genetics*. 1997; 146:1197–206. [PubMed: 9215920]
64. Weir BS, Cockerham CC. Estimating F-Statistics for the Analysis of Population Structure. *Evolution*. 1984; 38:1358–1370.
65. Chen H, Patterson N, Reich D. Population differentiation as a test for selective sweeps. *Genome research*. 2010; 20:393–402. [PubMed: 20086244]
66. Hofer T, Ray N, Wegmann D, Excoffier L. Large allele frequency differences between human continental groups are more likely to have occurred by drift during range expansions than by selection. *Annals of human genetics*. 2009; 73:95–108. [PubMed: 19040659]
67. Stephens M, Smith NJ, Donnelly P. A new statistical method for haplotype reconstruction from population data. *American journal of human genetics*. 2001; 68:978–89. [PubMed: 11254454]
68. Stephens M, Scheet P. Accounting for decay of linkage disequilibrium in haplotype inference and missing-data imputation. *American journal of human genetics*. 2005; 76:449–62. [PubMed: 15700229]
69. Climer S, Zhang W. Rearrangement Clustering: Pitfalls, remedies, and applications. *Journal of Machine Learning Research*. 2006; 7:919–943.
70. Cook, WJ. In *Pursuit of the Traveling Salesman: Mathematics at the Limits of Computation*. Princeton University Press; 2011. at <<http://press.princeton.edu/titles/9531.html>>

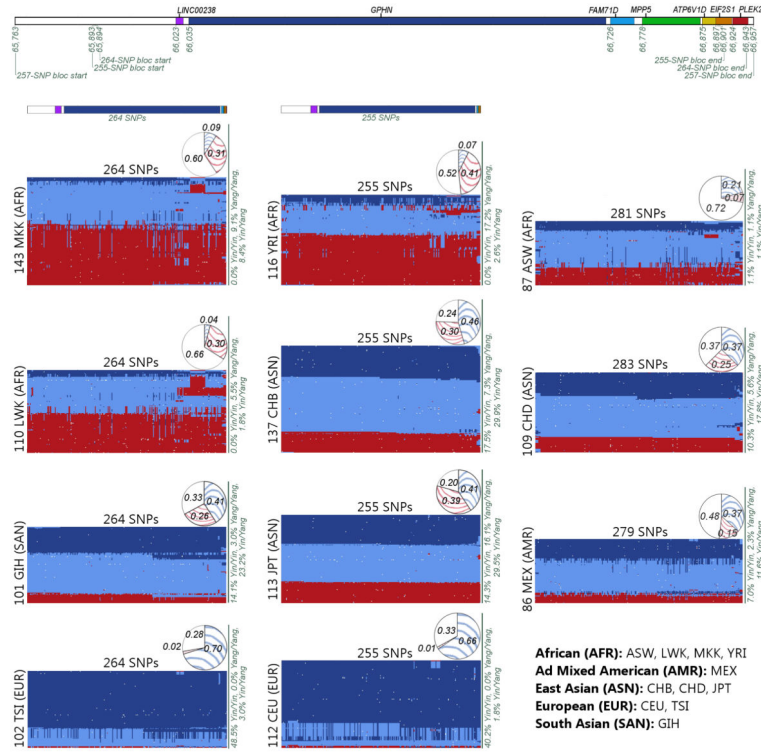


Figure 1. Yin-yang haplotypes

(Best viewed in color, high-resolution image available online as Supplementary Fig. 3.)

Upper panel: Yin-yang region with color-coded genes and positions in kb. *Lower panel:* Genotypes for the yin-yang haplotypes. The first column of matrices represents the ‘yin’ bloc identified in the analysis of the GIH, LWK, MKK, and TSI populations. The middle column represents the ‘yin’ bloc from the CEU, CHB, JPT, and YRI analysis. The last column represents the available yin-yang SNPs for each of the three additional HapMap populations: ASW, CHD, and MEX. For each matrix, each column represents a SNP and each row represents an individual (rows are rearranged to place similar individuals near each other). Color-coded bars at top of the first two columns indicate SNP positions by matching the gene color from the top panel. Dark blue indicates homozygote for SNP allele in the bloc, red for homozygote for alternate allele, light blue for heterozygote, and white for missing data. A solid dark blue horizontal line represents an individual that possesses two yin haplotypes and a solid red line represents a yang homozygote. Percentages of individuals that are homozygotes or heterozygotes for the yin and yang haplotypes are shown on the right side of each matrix. Yin (blue) and yang (red) haplotype frequencies are shown in pie chart above each matrix, with white indicating the percentage of haplotypes that are not 100% yin nor 100% yang.

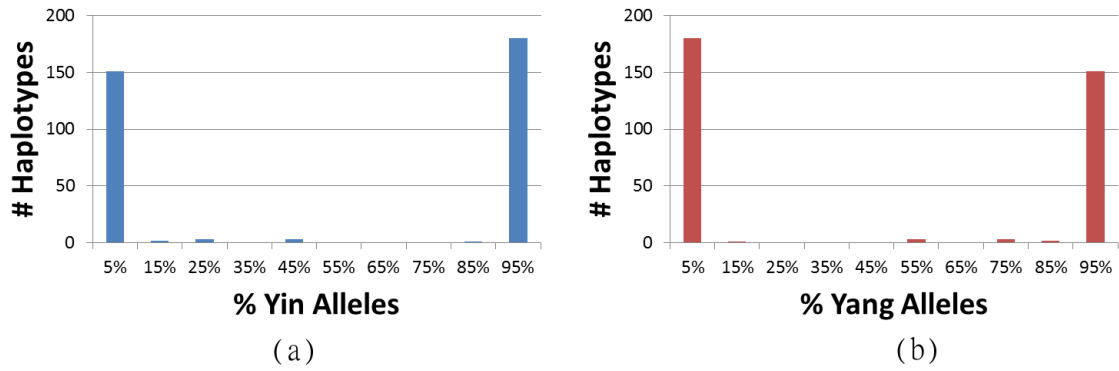


Figure 2. Haplotype compositions

The numbers of CHB and JPT phased chromosomes possessing various percentages of (a) yin and (b) yang SNP alleles are shown. Because these are biallelic SNPs, the plots are mirror images (e.g. the haplotypes representing the 0%–10% range in (a) are the same haplotypes representing the 90%–100% range in (b)). Only nine of the 340 phased chromosomes lie in the 10%–90% range for yin or yang alleles.

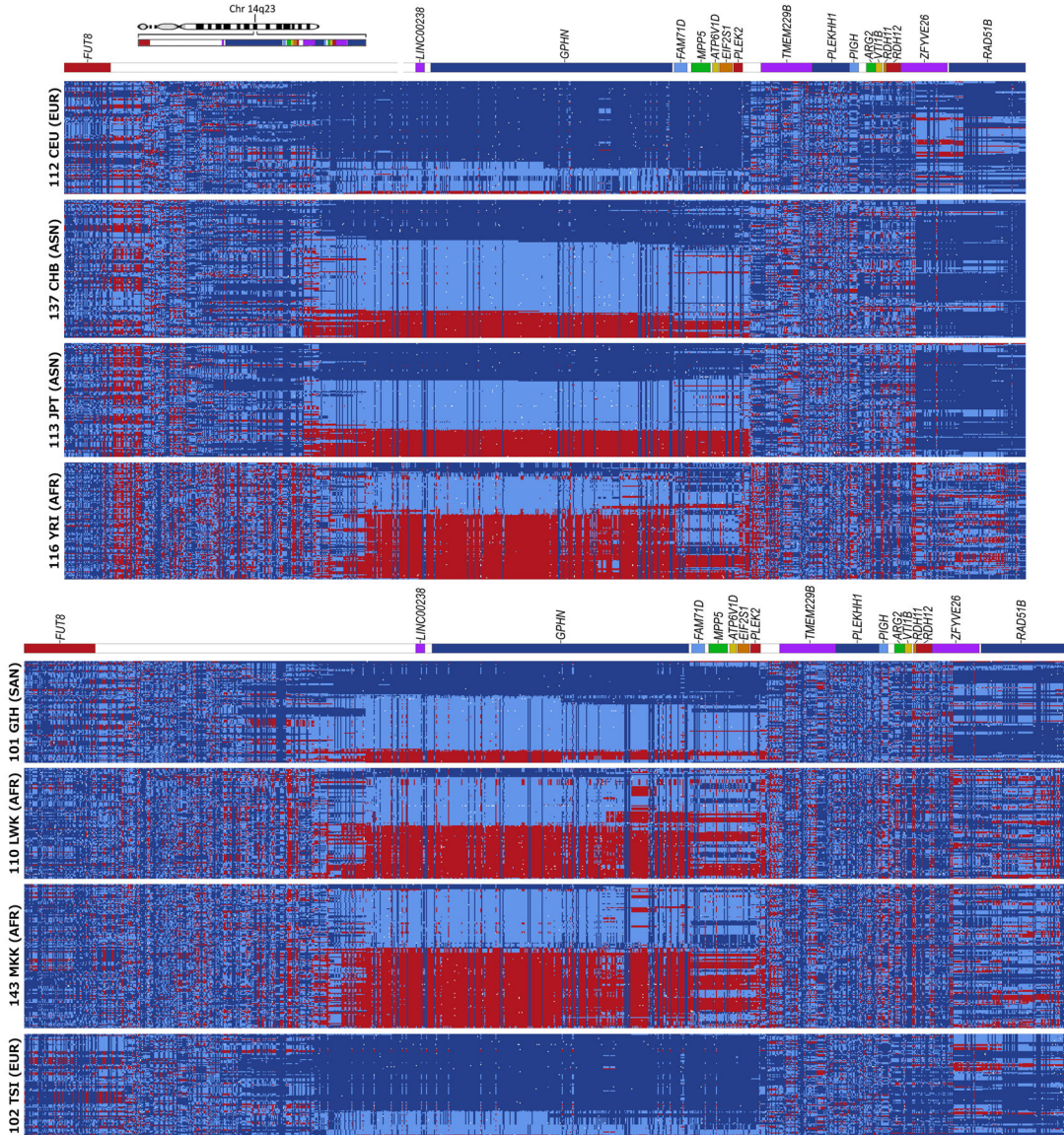


Figure 3. All SNPs within and surrounding yin-yang region
 The available HapMap SNPs within the yin-yang region (355 for the first four populations and 436 for the second four populations) are shown, along with 300 SNPs upstream and 300 downstream. (Best viewed in color; see Fig. 1 caption for details.) Homozygotes in the minor allele for the CEU or TSI populations are colored red, heterozygotes are light blue, homozygotes in the alternate allele are dark blue, and missing data are white. SNPs with low minor allele frequencies (MAF) appear as vertical columns that are predominantly either dark blue or red. Color bars above each group indicate the distributions of the SNPs across genomic regions. High resolution images are available online as Supplementary Figs. 4 – 11.

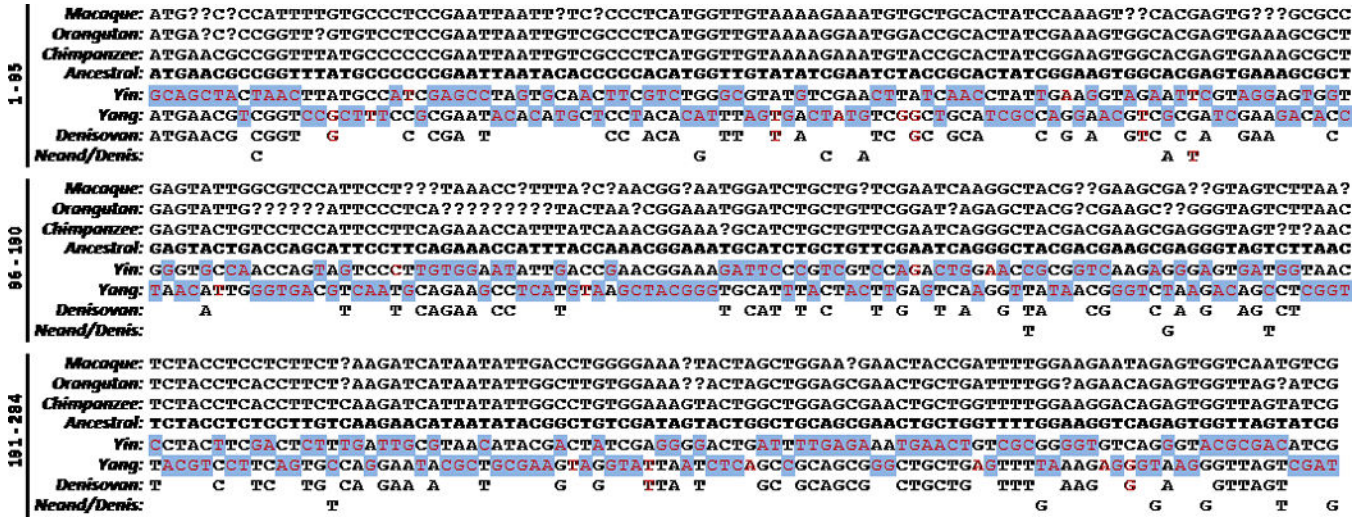


Figure 4. Haplotype comparisons

Reference alleles from the UCSC Genome Browser database are shown for macaque, orangutan, and chimpanzee. The ancestral alleles are supplied from NCBI’s dbSNP website (see Methods). Red font indicates variation from these ancestral alleles for the yin and yang haplotypes, as well as the haplotype drawn from the individual from the Denisovan fossil site²⁸ and the 15 SNP alleles representing several Neandertal and the Denisovan individual (Neand/Denis, see Methods). All 125 identified genotypes for the Denisovan haplotype are homozygous for the allele shown and this haplotype matches the yang haplotype at all 125 sites. However it is also highly similar to the ancestral haplotype, with 95.2% of the alleles matching. On the other hand, 11 of the 15 Neand/Denis alleles match the yin alleles while the remaining four match the ancestral, chimpanzee, orangutan, and macaque alleles. All but one of the Neand/Denis alleles match the ancestral haplotype and the derived allele matches the yin haplotype. Note that the variation that defines the yin and yang haplotypes (red font) is predominantly unique as 92.6% of the 284 alleles do not match the Neand/Denis, Denisovan, chimpanzee, or orangutan alleles (blue shading). Both the yin and yang haplotypes are dissimilar from the ancestral haplotype, having only about half of the SNP alleles in common, yet these two haplotypes are highly conserved across modern human populations. See Supplementary Dataset 1 for additional information.

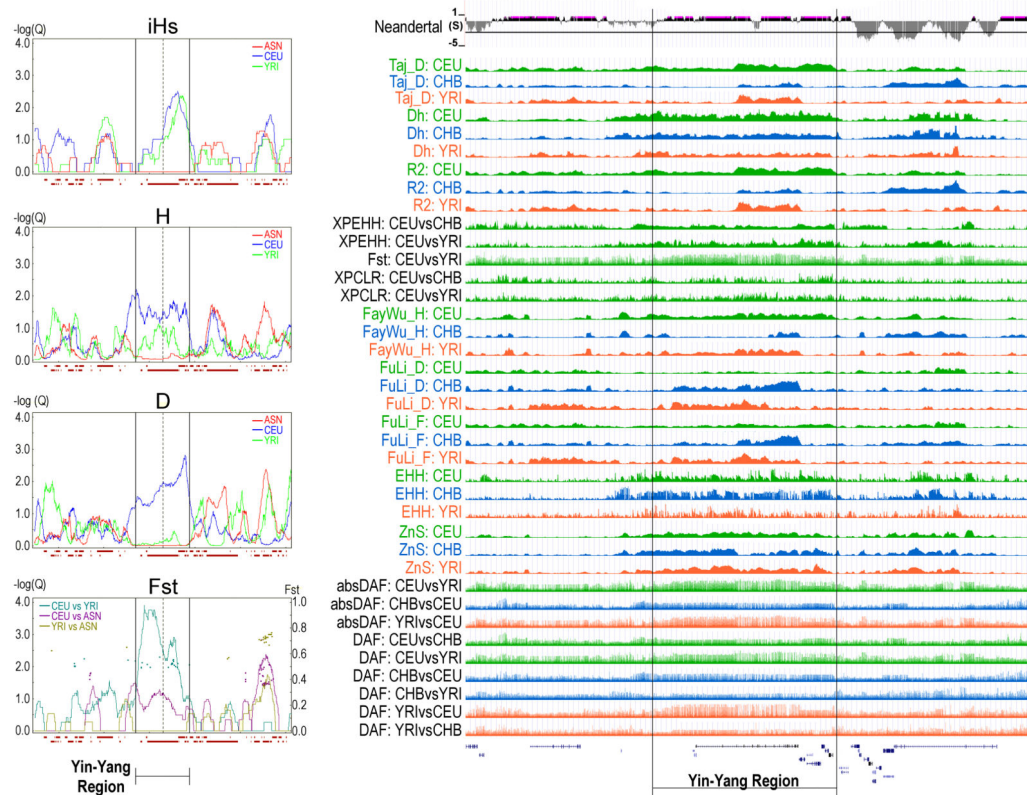


Figure 5. Statistical tests for selection

Results from Haplotter²² (left) and The 1000 Genomes Selection Browser 1.0²⁹ (right). See Methods for descriptions of tests. Left panel displays results for selection tests over a 5-Mb window centered on *gephyrin* for four HapMap populations, CEU, YRI, CHB, and JPT. The CHB and JPT are combined and labeled as ‘ASN’ in the plots. These plots include Voight *et al.*’s *iHs*, Fay and Wu’s *H*, Tajima’s *D*, and F_{ST} . Topmost plot on right represents a selective sweep scan on Neandertal vs. human polymorphisms using *Z*-score +-variance (Neandertal), followed by rank scores for 13 tests for selection. The rank scores were computed using an outlier approach based on sorted genome-wide scores and peaks represent regions under positive selection²⁹. Included are Tajima’s *D* (Taj_D), Nei’s *Dh* (Dh), Ramos-Onsins and Rozas’ R^2 (R2), Sabeti *et al.*’s XP-EHH* (XPEHH), Weir and Cockerham’s pairwise F_{ST} (Fst), Chen *et al.*’s XP-CLR (XPCLR), Fay and Wu’s *H* (FayWu_H), Fu and Li’s *D* (FuLi_D), Fu and Li’s *F* (FuLi_F), Sabeti *et al.*’s EHH_average* (EHH), Kelly’s ZnS (ZnS), Hofer *et al.*’s absolute DAF (absDAF), and Hofer *et al.*’s standard DAF (DAF). Asterisk indicates the method was modified by Pybus *et al.*²⁹ Populations are indicated. This image includes modified screen shots from <http://hsb.upf.edu/> and <http://haplotter.uchicago.edu/>.

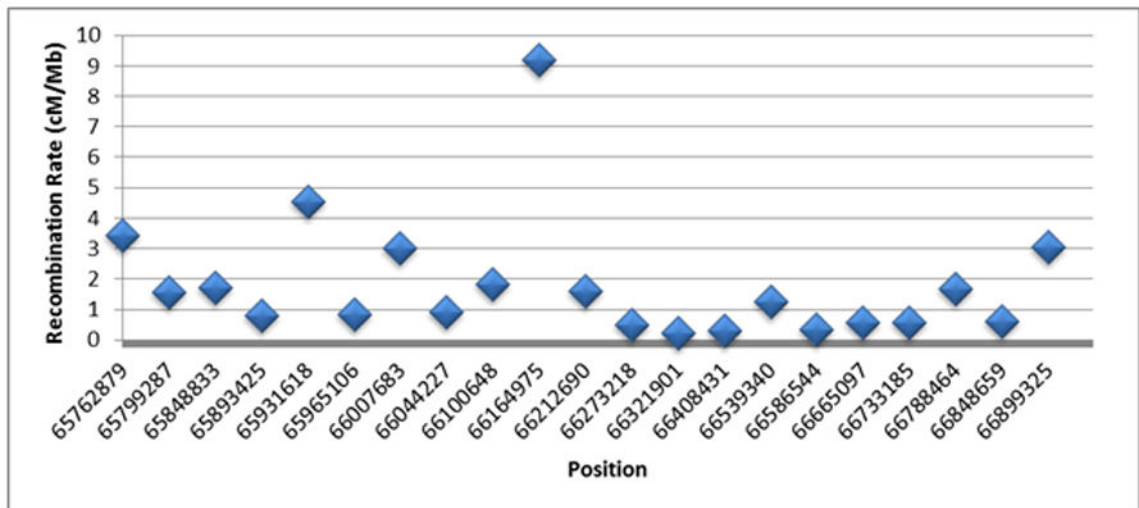


Figure 6. Recombination rates

Shown are the recombination rates provided by the HapMap Consortium for the yin-yang region. These rates were computed using the CEU, CHB, JPT, and YRI population data. All of these positions are within the yin-yang region and the recombination rate of 9.2 cM Mb⁻¹ was estimated at rs10133120 in the 5' end of *gephyrin*.