

28. J. Savarino, M. Legrand, *J. Geophys. Res.* **103**, 8267 (1998).
29. P. D. Jones, M. E. Mann, *Rev. Geophys.* **42**, 2003RG000143 (2004).
30. A. Moberg, D. M. Sonechkin, K. Holmgren, N. M. Datsenko, W. Karlen, *Nature* **433**, 613 (2005).
31. M. O. Andreae, P. Merlet, *Global Biogeochem. Cycles* **15**, 955 (2001).
32. C. McEvedy, R. Jones, *Atlas of World Population History* (Penguin, London, 1978), pp. 368.
33. W. Denevan, *The Native Population of the Americas in 1492* (Univ. of Wisconsin Press, Madison, WI, 1992), pp. 386.
34. B. Glaser, L. Haumaier, G. Guggenberger, W. Zech, *Naturwissenschaften* **88**, 37 (2001).
35. C. MacFarling Meure, thesis, University of Melbourne (2004).
36. We thank the staff of the Australian Antarctic Program, especially Casey Station, for field support; A. Smith for firn-air sampling assistance; the Bureau of Meteorology (Australia) for Cape Grim archive-air collection assistance; R. Francey, P. Steele, C. Allison, and S. Coram at CSIRO for logistical and technical help; and especially B. Ruddiman and B. Allan for valuable discussions. Supported by NSF (grant no. OPP0087357); NOAA/Climate Modeling and Diagnostics Laboratory; NIWA, New Zea-

land (Foundation for Research Science and Technology grant no. C01X0204); and the Australian Government's Antarctic Climate and Ecosystems Cooperative Research Centre and CSIRO Atmospheric Research.

#### Supporting Online Material

www.sciencemag.org/cgi/content/full/309/5741/1714/DC1  
Materials and Methods  
References and Notes

23 May 2005; accepted 28 July 2005  
10.1126/science.1115193

## Microcephalin, a Gene Regulating Brain Size, Continues to Evolve Adaptively in Humans

Patrick D. Evans,<sup>1,2</sup> Sandra L. Gilbert,<sup>1</sup> Nitzan Mekel-Bobrov,<sup>1,2</sup> Eric J. Vallender,<sup>1,2</sup> Jeffrey R. Anderson,<sup>1</sup> Leila M. Vaez-Azizi,<sup>1</sup> Sarah A. Tishkoff,<sup>4</sup> Richard R. Hudson,<sup>3</sup> Bruce T. Lahn<sup>1\*</sup>

The gene *Microcephalin* (*MCPH1*) regulates brain size and has evolved under strong positive selection in the human evolutionary lineage. We show that one genetic variant of *Microcephalin* in modern humans, which arose ~37,000 years ago, increased in frequency too rapidly to be compatible with neutral drift. This indicates that it has spread under strong positive selection, although the exact nature of the selection is unknown. The finding that an important brain gene has continued to evolve adaptively in anatomically modern humans suggests the ongoing evolutionary plasticity of the human brain. It also makes *Microcephalin* an attractive candidate locus for studying the genetics of human variation in brain-related phenotypes.

The most distinct trait of *Homo sapiens* is the exceptional size and complexity of the brain (1, 2). Several recent studies have linked specific genes to the evolution of the human brain (3–12). One of these is *Microcephalin* (7, 8); mutations in this gene cause primary microcephaly [MCPH; Online Mendelian Inheritance in Man (OMIM) accession 251200] (13, 14). MCPH is defined clinically as severe reductions in brain size coupled with mental retardation, but remarkably, an overall retention of normal brain structure and a lack of overt abnormalities outside of the nervous system (15–17). This led to the notion that the brains of MCPH patients function normally for their size and that genes underlying MCPH are specific developmental regulators of brain size (15–17).

*Microcephalin* is one of six known loci, named *MCPH1* through *MCPH6*, for which recessive mutations lead to MCPH (14, 18–23). For four of these, the underlying genes have been identified as *Microcephalin* (*MCPH1*), *CDK5RAP2* (*MCPH3*), *ASPM* (*MCPH5*), and

*CENPJ* (*MCPH6*) (14, 21, 23). Patients with loss-of-function mutations in *Microcephalin* have cranial capacities about 4 SD below the mean at birth. As adults, their typical brain size is around 400 cm<sup>3</sup> (whereas the normal range is 1200 to 1600 cm<sup>3</sup>), and the cerebral cortex is especially small (13, 14). *Microcephalin* is suggested to control the proliferation and/or differentiation of neuroblasts during neurogenesis. This postulate was consistent with several observations. First, mouse *Microcephalin* is expressed prominently in the proliferative zones of the embryonic brain (14). Second, the *Microcephalin* protein contains several copies of the BRCT domain that is found in cell cycle regulators, such as *BRCA1* (14, 24). Finally, cell culture studies indeed suggested a role of *Microcephalin* in regulating cell cycle (25–27).

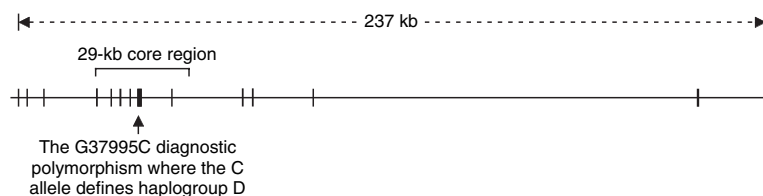
The finding that *Microcephalin* is a critical regulator of brain size spurred the hypothesis

that it might have played a role in brain evolution (16, 28). Consistent with this hypothesis, phylogenetic analysis of *Microcephalin* revealed signatures of strong positive selection in the lineage leading to humans (7, 8). Here, we examine the possibility that positive selection has continued to operate on this gene after the emergence of anatomically modern humans.

The human *Microcephalin* locus has 14 exons spanning about 236 kb on chromosome 8p23 (14) (Fig. 1). We previously sequenced all the exons in 27 humans (8). When re-analyzing the data, we noticed that one haplotype had a much higher frequency than the other haplotypes. Additionally, this haplotype differed consistently from the others at position 37995 of the genomic sequence (counting from the start codon) or position 940 of the open reading frame. This polymorphism falls in exon 8 and changes amino acid residue 314 from an ancestral aspartate to a histidine. (This polymorphism is described as G37995C with G denoting the ancestral allele.)

To investigate whether positive selection has acted on the high-frequency haplotype, we resequenced 23.4 kb of a 29-kb region centered around the G37995C polymorphism (Fig. 1). Sequencing was performed on a panel of 89 individuals from the Coriell Institute, which broadly represents human diversity (see SOM). To assign the ancestral state of polymorphisms, we also sequenced the common chimpanzee. Several GC-rich segments were not sequenced because of technical difficulties. The resulting sequence data contained 220 polymorphic sites, including 213 single-nucleotide polymorphisms (SNPs) and 7 insertion/deletion polymorphisms (indels) (table S1).

Haplotypes were inferred using the PHASE 2.1 program (29, 30). A total of 86 haplotypes



**Fig. 1.** Genomic structure of the human *Microcephalin* gene. The region sequenced in the 89-individual Coriell panel is bracketed.

<sup>1</sup>Howard Hughes Medical Institute, Department of Human Genetics, <sup>2</sup>Committee on Genetics, <sup>3</sup>Department of Ecology and Evolution, University of Chicago, Chicago, IL 60637, USA. <sup>4</sup>Department of Biology, University of Maryland, College Park, MD 20742, USA.

\*To whom correspondence should be addressed. E-mail: blahn@bsd.uchicago.edu

were identified along with their frequencies (Fig. 2 and table S2). One haplotype, denoted 49, had a much higher frequency than the other haplotypes. It had the derived C allele at the G37995C SNP site and corresponded to the high-frequency haplotype in the aforementioned exon-only polymorphism survey (8). In the Coriell panel, haplotype 49 had a frequency of 33% (59 out of 178 chromosomes) and is found in all the populations sampled in the panel. The remaining 85 haplotypes varied in frequency from 0.6 to 6.2% (1 to 11 chromosomes).

Positive selection on an allele can increase the frequency of the haplotype bearing the allele while maintaining extended linkage disequilibrium (LD) around that allele (31–36). Our data on haplotype 49 are consistent with these signatures of selection. We formally tested the statistical significance of positive selection using the previously established coalescent model (37, 38). Given the slight uncertainty in haplotype inference, we considered only the 18 individuals in the Coriell panel who are homozygous for haplotype 49 (table S1).

By simulation, we calculated the probability of obtaining 18 or more individuals (out of 89) who are homozygous for a single haplotype across a region of 220 segregating sites under neutral evolution. Here, recombination and gene conversion rates were set to values previously established for the *Microcephalin* locus (39), and a demographic model with a severe bottleneck followed by exponential growth was assumed (see SOM). Prior studies have shown that the bottleneck specified here is likely to be much more stringent than that associated with the real demographic history of human populations (40, 41); thus, the test is conservative (38). Under these parameters, the probability of obtaining 18 homozygotes out of 89 is highly significant ( $P = 0$  based on 5,000,000 replicates).

We then tested several additional demographic models, including (i) constant size, (ii) very ancient expansion, (iii) very recent expansion, (iv) repeated severe bottlenecks with subsequent expansion, and (v) population structure with between two and five subpopulations (see SOM). All produced exceedingly significant results. Even though the exact demographic history of humans is yet to be defined, our tests are highly significant under a broad range of demographic scenarios, which furthers the argument that the statistical significance is unlikely to be altered by reasonable variations in the supposed human demography. We also tested the significance of the inferred haplotype data (i.e., the significance of having 59 copies of haplotype 49 among 178 chromosomes), which similarly produced highly significant results. These data strongly suggest that haplotype 49 was driven to high frequency by positive selection. However, our data do not address whether the positive selection is

frequency-dependent selection, heterozygote advantage, or simple additive positive selection.

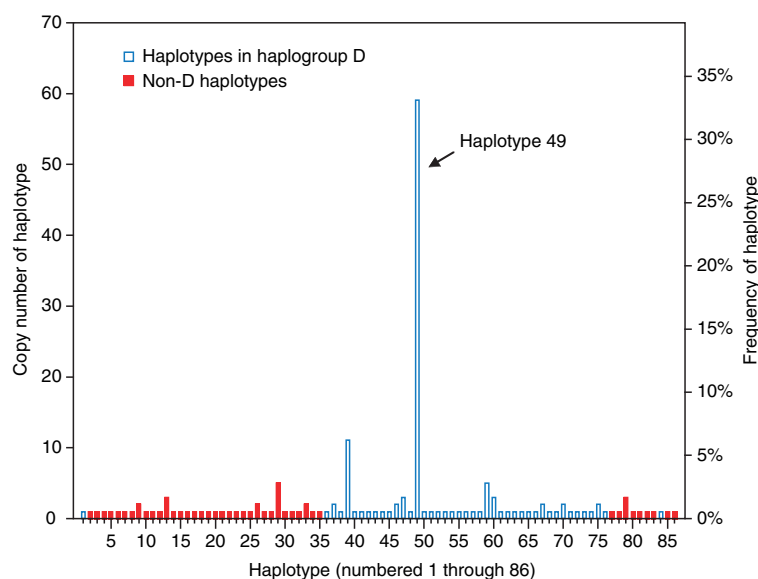
Using the G37995C polymorphism as a diagnostic site, we divided all the haplotypes into two groups: those that carry the derived C allele and those that carry the ancestral G allele. We designated the former group as haplogroup D (where D stands for “derived”). It includes 43 haplotypes that together have a 70% frequency in the Coriell panel, and haplotype 49 is the predominant member (table S2). Although the derived C allele at the G37995C site only provides an operational definition for haplogroup D, several observations make evident that haplogroup D is systematically different from the non-D haplotypes. First, this haplogroup consists exclusively of haplotype 49 or its minor variants, whereas non-D haplotypes show much greater sequence divergence from haplogroup D chromosomes. This greater divergence is because haplogroup D and non-D haplotypes have multiple fixed differences relative to each other in addition to G37995C (table S2). The only exceptions are a few recombinant haplotypes between D and non-D chromosomes (discussed below). Second, for sites that are polymorphic within haplogroup D chromosomes (excluding recombinants between D and non-D chromosomes), the non-D chromosomes are invariably monomorphic for the ancestral alleles. These data indicate that haplogroup D constitutes a genealogical clade of closely related haplotypes that is altogether separate from the more distantly related non-D haplotypes (again, excluding recombinants between D and non-D chromosomes, which represent mixed genealogies).

Collectively, the above observations support an evolutionary scenario with two aspects.

First, haplotype 49 swept from a single copy to high frequency in a short period of time. Second, during the sweep, minor variants of haplotype 49 emerged through rare mutations and recombinations. These variants, together with haplotype 49, make up haplogroup D. Haplotype 49 evidently represents the most recent common ancestor (MRCA) of haplogroup D, because it consistently has the ancestral allele for the sites polymorphic within haplogroup D.

We next estimated the coalescence age (i.e., time to MRCA) of haplogroup D chromosomes in the Coriell panel. We used the average number of mutations from the MRCA of a haplogroup clade to its descendant lineages as a molecular clock for estimating the age of the clade (42, 43). This approach is known to be unbiased by demographic history (42). The age of haplogroup D was found to be ~37,000 years, with a 95% confidence interval of 14,000 to 60,000 years. In comparison, the coalescence age of all the chromosomes in the Coriell panel is about 1,700,000 years. The emergence of anatomically modern humans has been estimated to be 200,000 years before present (44). Haplogroup D is obviously much younger, which indicates that positive selection was at work in a period considerably postdating the emergence of anatomically modern humans in Africa. We note that the age of haplogroup D coincides with the introduction of anatomically modern humans into Europe about 40,000 years ago, as well as the dramatic shift in the archeological record indicative of modern human behavior, such as art and the use of symbolism (i.e., the “Upper Paleolithic revolution”) (45).

If haplogroup D indeed experienced a recent selective sweep, it should show low poly-



**Fig. 2.** Frequencies of 86 inferred *Microcephalin* haplotypes in the 89-individual Coriell panel. Haplotypes in haplogroup D are indicated by blue-edged bars; non-D haplotypes are indicated by solid red bars.

morphism and an excess of rare alleles (46). To confirm this, we calculated nucleotide diversity ( $\pi$ ) and Tajima's  $D$  for the 47 individuals who are homozygous for haplogroup D chromosomes, and we compared these values to those of the non-D chromosomes. The  $\pi$  value of the D chromosomes is lower, by a factor of 12, than that of the non-D chromosomes (0.000077 and 0.00092, respectively), even though the D chromosomes represent about 70% of the chromosomes in the panel. Tajima's  $D$ , which is a summary statistic for the frequency spectrum of alleles, is  $-2.3$  for haplogroup D (whereas it is  $-1.2$  for the non-D chromosomes). This strongly negative Tajima's  $D$  indicates a starlike genealogy for haplogroup D chromosomes (47). Thus, both summary statistics contrast sharply between D and non-D chromosomes and are consistent with the recent age and rapid expansion of haplogroup D. We note that these calculations do not provide a statistically stringent test of positive selection, because they are done on subsets of the genealogy. Nevertheless, they do

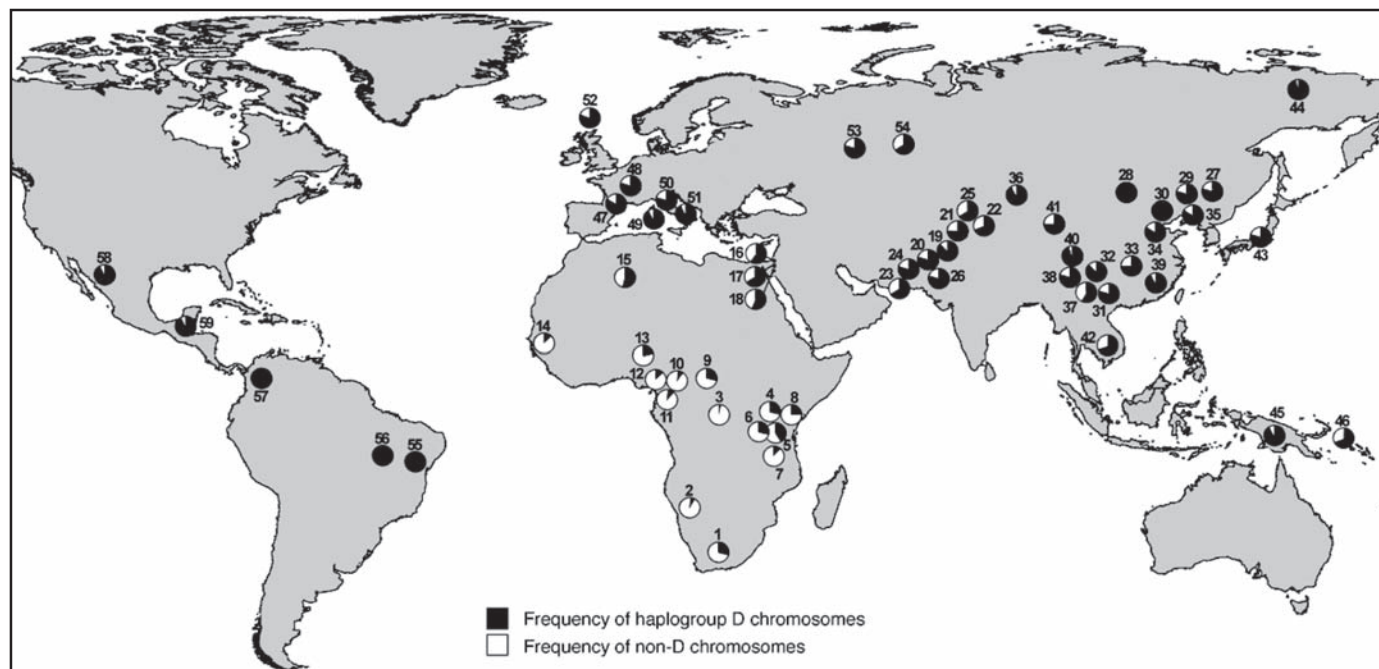
reveal qualitative signatures of positive selection that further corroborate the more stringent statistical tests described earlier.

Another sign of a positive selective sweep is extended LD around the selected allele. This is apparent in the region of *Microcephalin* investigated here, where haplogroup D chromosomes show near-complete LD across the entire region. The only exceptions are haplotypes 1, 68, and 84 (each found in a single copy in the Coriell panel), which are recombinants between D and non-D chromosomes as evidenced by recombination tracts (table S2). The remaining 121 copies of haplogroup D chromosomes show no evidence of recombination. By comparison, the non-D chromosomes do not display any significant LD across the region.

To probe the extent of LD beyond the 29-kb core region, we sequenced the Coriell panel for two segments of about 3 kb each, situated at the beginning and end of the gene separated from each other by about 235 kb. In these flanking regions, there is clear evidence

of LD decay from the core region, which supports the idea that selection has most likely operated on a site (or sites) around the core region. Our present data cannot resolve the exact site(s) of selection, and the G37995C nonsynonymous SNP used to define haplogroup D is just a candidate.

To obtain a more detailed frequency distribution of haplogroup D across the globe, we analyzed a much larger human population panel containing 1184 globally diverse individuals. We genotyped the diagnostic G37995C SNP in this panel to infer the frequency of haplogroup D chromosomes (Fig. 3). Geographic variation was observed, with sub-Saharan populations generally having lower frequencies than others. The statistic for genetic differentiation,  $F_{ST}$ , is 0.48 between sub-Saharan and others, which indicates strong differentiation (48) and is significantly higher than the genome average of 0.12 ( $P < 0.03$  based on previously established genome-wide  $F_{ST}$  distribution) (49). Such population differentiation may reflect a Eurasian origin of haplogroup D, local adaptation, and/or



**Fig. 3.** Global frequencies of *Microcephalin* haplogroup D chromosomes (defined as having the derived C allele at the G37995C diagnostic SNP) in a panel of 1184 individuals. For each population, the country of origin, number of individuals sampled, and frequency of haplogroup D chromosomes are given (in parentheses) as follows: 1, Southeastern and Southwestern Bantu (South Africa, 8, 31.3%); 2, San (Namibia, 7, 7.1%); 3, Mbuti Pygmy (Democratic Republic of Congo, 15, 3.3%); 4, Masai (Tanzania, 27, 29.6%); 5, Sandawe (Tanzania, 32, 39.1%); 6, Burunge (Tanzania, 28, 30.4%); 7, Turu (Tanzania, 23, 15.2%); 8, Northeastern Bantu (Kenya, 12, 25%); 9, Biaka Pygmy (Central African Republic, 32, 26.6%); 10, Zime (Cameroon, 23, 8.7%); 11, Bakola Pygmy (Cameroon, 24, 10.4%); 12, Bamoun (Cameroon, 28, 17.9%); 13, Yoruba (Nigeria, 25, 24%); 14, Mandenka (Senegal, 24, 16.7%); 15, Mozabite [Algeria (Mzab region), 29, 53.5%]; 16, Druze [Israel (Carmel region), 44, 60.2%]; 17, Palestinian [Israel (Central), 40, 63.8%]; 18, Bedouin [Israel (Negev region), 44, 54.6%]; 19, Hazara (Pakistan, 20, 85%); 20, Balochi (Pakistan, 23, 78.3%); 21, Pathan (Pakistan, 23, 76.1%); 22, Burusho (Pakistan, 25, 66%); 23, Makrani (Pakistan, 24,

62.5%); 24, Brahui (Pakistan, 25, 78%); 25, Kalash (Pakistan, 24, 62.5%); 26, Sindhi (Pakistan, 25, 78%); 27, Hezhen (China, 9, 77.8%); 28, Mongola (China, 10, 100%); 29, Daur (China, 10, 85%); 30, Orogen (China, 10, 100%); 31, Miaozi (China, 9, 77.8%); 32, Yizu (China, 10, 85%); 33, Tujia (China, 10, 75%); 34, Han (China, 41, 82.9%); 35, Xibo (China, 9, 83.3%); 36, Uygur (China, 10, 90%); 37, Dai (China, 9, 55.6%); 38, Lahu (China, 10, 85%); 39, She (China, 9, 88.9%); 40, Naxi (China, 10, 95%); 41, Tu (China, 10, 75%); 42, Cambodian (Cambodia, 11, 72.7%); 43, Japanese (Japan, 27, 77.8%); 44, Yakut [Russia (Siberia region), 25, 98%]; 45, Papuan (New Guinea, 17, 91.2%); 46, NAN Melanesian (Bougainville, 18, 72.2%); 47, French Basque (France, 24, 83.3%); 48, French (France, 28, 78.6%); 49, Sardinian (Italy, 26, 90.4%); 50, North Italian [Italy (Bergamo region), 13, 76.9%]; 51, Tuscan (Italy, 8, 87.5%); 52, Orcadian (Orkney Islands, 16, 81.3%); 53, Russian (Russia, 24, 79.2%); 54, Adygei [Russia (Caucasus region), 15, 63.3%]; 55, Karitiana (Brazil, 21, 100%); 56, Surui (Brazil, 20, 100%); 57, Colombian (Colombia, 11, 100%); 58, Pima (Mexico, 25, 92%); 59, Maya (Mexico, 25, 92%).



demographic factors such a bottleneck associated with human migration out of Africa 50,000 to 100,000 years ago.

Previous studies have shown that *Microcephalin* is a specific regulator of brain size (13, 14) and that this gene has evolved under strong positive selection in the primate lineage leading to *Homo sapiens* (7, 8). Here, we present compelling evidence that *Microcephalin* has continued its trend of adaptive evolution beyond the emergence of anatomically modern humans. The specific function of *Microcephalin* in brain development makes it likely that selection has operated on the brain. Yet, it remains formally possible that an unrecognized function of *Microcephalin* outside of the brain is actually the substrate of selection. If selection indeed acted on a brain-related phenotype, there could be several possibilities, including brain size, cognition, personality, motor control, or susceptibility to neurological and/or psychiatric diseases. We hypothesize that D and non-D haplotypes have different effects on the proliferation of neural progenitor cells, which in turn leads to different phenotypic outcomes of the brain visible to selection.

References and Notes

1. J. N. Spuhler, *The Evolution of Man's Capacity for Culture* (Wayne State Univ. Press, Detroit, MI, 1959).
2. J. H. Jerison, *Evolution of the Brain and Intelligence* (Academic Press, New York, 1973).
3. W. Enard et al., *Nature* **418**, 869 (2002).

4. J. Zhang, *Genetics* **165**, 2063 (2003).
5. P. D. Evans et al., *Hum. Mol. Genet.* **13**, 489 (2004).
6. N. Kouprina et al., *PLoS Biol.* **2**, E126 (2004).
7. Y. Q. Wang, B. Su, *Hum. Mol. Genet.* **13**, 1131 (2004).
8. P. D. Evans, J. R. Anderson, E. J. Vallender, S. S. Choi, B. T. Lahn, *Hum. Mol. Genet.* **13**, 1139 (2004).
9. R. J. Ferland et al., *Nat. Genet.* **36**, 1008 (2004).
10. H. H. Stedman et al., *Nature* **428**, 415 (2004).
11. F. Burki, H. Kaessmann, *Nat. Genet.* **36**, 1061 (2004).
12. S. Dorus et al., *Cell* **119**, 1027 (2004).
13. A. P. Jackson et al., *Am. J. Hum. Genet.* **63**, 541 (1998).
14. A. P. Jackson et al., *Am. J. Hum. Genet.* **71**, 136 (2002).
15. W. B. Dobyns, *Am. J. Hum. Genet.* **112**, 315 (2002).
16. G. H. Mochida, C. A. Walsh, *Curr. Opin. Neurol.* **14**, 151 (2001).
17. C. G. Woods, J. Bond, W. Enard, *Am. J. Hum. Genet.* **76**, 717 (2005).
18. E. Roberts et al., *Eur. J. Hum. Genet.* **7**, 815 (1999).
19. L. Moynihan et al., *Am. J. Hum. Genet.* **66**, 724 (2000).
20. C. R. Jamieson, C. Govaerts, M. J. Abramowicz, *Am. J. Hum. Genet.* **65**, 1465 (1999).
21. J. Bond et al., *Nat. Genet.* **32**, 316 (2002).
22. G. F. Leal et al., *J. Med. Genet.* **40**, 540 (2003).
23. J. Bond et al., *Nat. Genet.* **37**, 353 (2005).
24. T. Huyton, P. A. Bates, X. Zhang, M. J. Sternberg, P. S. Freemont, *Mutat. Res.* **460**, 319 (2000).
25. S. Y. Lin, S. J. Elledge, *Cell* **113**, 881 (2003).
26. X. Xu, J. Lee, D. F. Stern, *J. Biol. Chem.* **279**, 34091 (2004).
27. M. Trimbom et al., *Am. J. Hum. Genet.* **75**, 261 (2004).
28. S. L. Gilbert, W. B. Dobyns, B. T. Lahn, *Nat. Rev. Genet.* **6**, 581 (2005).
29. M. Stephens, N. J. Smith, P. Donnelly, *Am. J. Hum. Genet.* **68**, 978 (2001).
30. M. Stephens, P. Donnelly, *Am. J. Hum. Genet.* **73**, 1162 (2003).
31. S. A. Tishkoff et al., *Science* **293**, 455 (2001).
32. P. C. Sabeti et al., *Nature* **419**, 832 (2002).
33. E. Wang et al., *Am. J. Hum. Genet.* **74**, 931 (2004).
34. T. Bersaglieri et al., *Am. J. Hum. Genet.* **74**, 1111 (2004).
35. E. E. Thompson et al., *Am. J. Hum. Genet.* **75**, 1059 (2004).
36. H. Stefansson et al., *Nat. Genet.* **37**, 129 (2005).

37. R. R. Hudson, *Oxf. Surv. Evol. Biol.* **7**, 1 (1990).
38. R. R. Hudson, *Bioinformatics* **18**, 337 (2002).
39. A. Kong et al., *Nat. Genet.* **31**, 241 (2002).
40. E. Zietkiewicz et al., *J. Mol. Evol.* **47**, 146 (1998).
41. H. Harpending, A. Rogers, *Annu. Rev. Genomics Hum. Genet.* **1**, 361 (2000).
42. H. Tang, D. O. Siegmund, P. Shen, P. J. Oefner, M. W. Feldman, *Genetics* **161**, 447 (2002).
43. R. Thomson, J. K. Pritchard, P. Shen, P. J. Oefner, M. W. Feldman, *Proc. Natl. Acad. Sci. U.S.A.* **97**, 7360 (2000).
44. I. McDougall, F. H. Brown, J. G. Fleagle, *Nature* **433**, 733 (2005).
45. R. G. Klein, *The Human Career: Human Biological and Cultural Origins* (Univ. of Chicago Press, Chicago, 1999).
46. M. Bamshad, S. P. Wooding, *Nat. Rev. Genet.* **4**, 99 (2003).
47. F. Tajima, *Genetics* **123**, 585 (1989).
48. S. Wright, *Evolution and the Genetics of Populations* (Univ. of Chicago Press, Chicago, 1978).
49. J. M. Akey, G. Zhang, K. Zhang, L. Jin, M. D. Shriver, *Genome Res.* **12**, 1805 (2002).
50. We thank the Coriell Institute for Medical Research, the Centre d'Etude du Polymorphisme Humain (CEPH), and A. Froment for human DNA samples. We thank H. M. Cann, S. Dorus, E. E. Eichler, N. M. Pearson, A. Di Rienzo, M. Kreitman, and J. K. Pritchard for technical support and/or helpful discussions. Supported in part by the Searle Scholarship and the Burroughs Wellcome Career Award (to B.T.L.), and David and Lucile Packard Career Award, the Burroughs Wellcome Career Award, and NSF grant BCS-0196183 (to S.A.T.).

Supporting Online Material

www.sciencemag.org/cgi/content/full/309/5741/1717/DC1  
 Materials and Methods  
 Tables S1 and S2  
 References and Notes

18 April 2005; accepted 14 June 2005  
 10.1126/science.1113722

# Ongoing Adaptive Evolution of *ASPM*, a Brain Size Determinant in *Homo sapiens*

Nitzan Mekel-Bobrov,<sup>1,2</sup> Sandra L. Gilbert,<sup>1</sup> Patrick D. Evans,<sup>1,2</sup> Eric J. Vallender,<sup>1,2</sup> Jeffrey R. Anderson,<sup>1</sup> Richard R. Hudson,<sup>3</sup> Sarah A. Tishkoff,<sup>4</sup> Bruce T. Lahn<sup>1\*</sup>

The gene *ASPM* (*abnormal spindle-like microcephaly associated*) is a specific regulator of brain size, and its evolution in the lineage leading to *Homo sapiens* was driven by strong positive selection. Here, we show that one genetic variant of *ASPM* in humans arose merely about 5800 years ago and has since swept to high frequency under strong positive selection. These findings, especially the remarkably young age of the positively selected variant, suggest that the human brain is still undergoing rapid adaptive evolution.

Homozygous null mutations of *ASPM* cause primary microcephaly, a condition characterized by severely reduced brain size with otherwise normal neuroarchitecture (1). Studies

have suggested that *ASPM* may regulate neural stem cell proliferation and/or differentiation during brain development, possibly by mediating spindle assembly during cell division (1, 2). Phylogenetic analysis of *ASPM* has revealed strong positive selection in the primate lineage leading to *Homo sapiens* (3–5), especially in the past 6 million years of hominid evolution in which *ASPM* acquired about one advantageous amino acid change every 350,000 years (4). These data argue that *ASPM*

may have contributed to human brain evolution (3–6). Here, we investigate whether positive selection has continued to operate on *ASPM* since the emergence of anatomically modern humans.

Human *ASPM* has 28 exons with a 10,434–base pair open reading frame (1) (fig. S1). We resequenced the entire 62.1-kb genomic region of *ASPM* in samples from 90 ethnically diverse individuals obtained through the Coriell Institute and from a common chimpanzee (7). This revealed 166 polymorphic sites (table S1). Using established methodology (7), we identified 106 haplotypes. One haplotype, numbered 63, had an unusually high frequency of 21%, whereas the other haplotypes ranged from 0.56% to 3.3% (fig. S2). Moreover, this haplotype differed consistently from the others at multiple polymorphic sites (save for a few rare haplotypes that are minor mutational or recombinational variants of haplotype 63, as discussed later) (table S2). Two of these polymorphic sites are nonsynonymous, both in exon 18, and are denoted A44871G and C45126A (numbers indicate genomic positions from the start codon, and letters at the beginning and end indicate ancestral and derived alleles, respectively). These two sites reside in a region of the open reading frame that was shown previously to have experienced par-

<sup>1</sup>Howard Hughes Medical Institute, Department of Human Genetics, <sup>2</sup>Committee on Genetics, <sup>3</sup>Department of Ecology and Evolution, University of Chicago, Chicago, IL 60637, USA. <sup>4</sup>Department of Biology, University of Maryland, College Park, MD 20742, USA.

\*To whom correspondence should be addressed. E-mail: blahn@bsd.uchicago.edu

**Science Supporting Online Material*****Microcephalin*, a Gene Regulating Brain Size, Continues to Evolve Adaptively in Humans**

Patrick D. Evans, Sandra L. Gilbert, Nitzan Mekel-Bobrov, Eric J. Vallender, Jeffrey R. Anderson, Leila M. Vaez-Azizi, Sarah A. Tishkoff, Richard R. Hudson, Bruce T. Lahn

DOI: 10.1126/science.1113722

**Materials and Methods***Sequence acquisition and preliminary analysis*

A panel of 89 human samples from the Coriell Institute that broadly represent worldwide populations was used for resequencing. It includes 9 sub-Saharan Africans (Coriell numbers: 17341–17349), 7 North Africans (17378–17384), 9 Iberians (17091–17097, 17099, 17100), 7 Basques (15883–15887, 16185, 16188), 9 Russians (13820a, 13838, 13852, 13876, 13877, 13911–13914), 9 Middle Easterners (17331–17340), 9 South Asians (17021–17024, 17026–17030), 8 Chinese (16654, 16688, 16689, 17014, 17015, 17017–17019), 1 Japanese (11587), 8 Southeast Asians (17081, 17083, 17085–17090), 6 Pacific Islander (17385–17388, 17390, 17391), and 7 Andeans (17301, 17302, 17306–17310). A common chimpanzee (*Pan troglodytes*) was also included in the sequencing. Double-stranded sequences in regions of interest were obtained by PCR amplification followed by sequencing of PCR products. Sequenced regions include 24750–26292, 26988–29992, 30561–32132, 32841–42938, 43006–44351, 45808–49406, 50123–50908, and 52305–53776 (the first base of the initiation codon of *Microcephalin* is defined as position 1). The core region used for haplotype analysis spans 29027 bases (24750–53776), of which 23416 bases were sequenced. Sequence chromatograms were aligned by the Sequencher software (Gene Codes Corporation, Ann Arbor, MI). Polymorphisms were detected by direct visual inspection of sequence chromatograms. The ancestral alleles of polymorphisms were called using the chimpanzee sequence as outgroup. Inference of haplotypes from the diploid sequence data was performed using the PHASE 2.1 software as described (S1, S2), which is available online at <http://www.stats.ox.ac.uk/mathgen/home.html>. Nucleotide diversity ( $\pi$ ) and Tajima's *D* were calculated using the program DnaSP 3.51, as described previously (S3). To avoid uncertainties of haplotype inference, the 47 individuals who are homozygous for haplogroup D chromosomes were used for the calculation of  $\pi$  and Tajima's *D* of this haplogroup. Inferred haplotypes were used to calculate  $\pi$  and Tajima's *D* for the non-D chromosomes. Recombinants between D and non-D chromosomes were excluded from the calculation.

*Genotyping*

Genotyping of the G37995C nonsynonymous polymorphism in *Microcephalin* was performed on a panel of 1184 human samples. This panel does not overlap with the Coriell panel described above. It consists of the HGDP CEPH diversity panel as described previously (S4), minus the following two sets of samples. One is a set of duplicated samples that needed to be removed, including HGDP00472, HGDP00452, HGDP00457, HGDP00980, HGDP00650, HGDP00583, HGDP00111, HGDP00220, HGDP00813, HGDP01233, HGDP00762, HGDP00770, HGDP00657, HGDP00658, HGDP00660, and HGDP01149. The other is a set of samples that failed to be

genotyped due to technical reasons (e.g., poor DNA quality), including HGDP01263, HGDP00633, HGDP00635, HGDP00636, HGDP00644, HGDP00579, HGDP00581, HGDP00584, HGDP00698, HGDP00700, HGDP00722, HGDP00723, HGDP00724, HGDP00725, HGDP00730, HGDP00731, HGDP00732, HGDP00734, HGDP00746, HGDP00076, HGDP00090, HGDP00109, HGDP00115, HGDP00122, HGDP00125, HGDP00141, HGDP00254, HGDP00281, HGDP00782, HGDP00783, HGDP01023, HGDP01193, HGDP01311, HGDP01334, HGDP00766, HGDP00768, HGDP00662, HGDP00520, HGDP00666, HGDP01077, HGDP01386, HGDP01402, HGDP00890, HGDP00707, HGDP00708, HGDP00995, HGDP00998, HGDP01010, and HGDP00841. The CEPH panel originally contained 1064 individuals, and had 999 individuals remaining after removing the above two sets of samples. Demographic information for the HGDP CEPH diversity panel is available online at <http://www.cephb.fr>. In addition, the panel contained 185 sub-Saharan African samples collected by S. A. Tishkoff and A. Froment (sample collection was approved by the Institutional Review Board at the University of Maryland). The samples included 23 Turu, 32 Sandawe, 28 Burunge, and 27 Masai individuals from Tanzania; they also included 24 Bakola Pygmy, 28 Bamoun, and 23 Zime individuals from Cameroon. To perform genotyping, a small region encompassing the G37995C polymorphism was amplified by PCR, followed by sequencing of the PCR product. Genotype was scored by visual inspection of the sequence chromatograms.  $F_{ST}$  was calculated as described previously (S5). The exact formulas are available on pages 143–155 of (S6).

### *Statistical analysis*

To test the statistical significance that the frequency of haplotype 49 departs from neutral expectation, we used a previously described simulation method based on the coalescent process as implemented in the ms software (S7, S8). First, the following parameters were specified: the number of chromosomes, the number of segregating sites, recombination rate, gene conversion rate, and demographic model. Recombination rate of the *Microcephalin* region was set at the locus-specific value of 1.9 cM/Mb as obtained in a previous genomewide survey (S9), and gene conversion rate was set to be the same as recombination rate with an average tract length of 100 bp. The gene conversion model was as previously described (S10), which assumes that the tract length is geometrically distributed. Nine demographic models were tested:

- 1) constant population with an effective size of  $10^4$ ,
- 2) an ancient population expansion from  $10^4$  at 5,000 generations ago exponentially to  $10^7$  today,
- 3) a recent population expansion from  $10^4$  at 1,000 generations ago exponentially to  $10^7$  today,
- 4) a severe bottle neck starting 5,000 generations ago that reduced the population from  $10^4$  instantly to  $10^3$  and lasted until 2,500 generations ago at which point the population started to expand exponentially to  $10^7$  today,
- 5) repeated bottlenecks for five successive rounds starting 7000 generations ago, each from  $10^4$  instantly to  $10^3$  for 500 generations followed by exponential recovery back to  $10^4$  over another 500 generations, except at the end of the fifth bottleneck 2500 generations ago which was followed by exponential growth to  $10^7$  today,
- 6) population structure where the initial 178 chromosomes were split equally into 2 different subpopulations under constant population size with 1 migration per generation, and
- 7 to 9) population structure where the initial 178 chromosomes were split equally into 3 to 5 different subpopulations with 1 migration per generation. Command lines in the ms program to input the above demographic models were as follows:

#### 1) Constant population size:

```
./ms 178 100000 -s 220 -r 11.6104 29027 -c 1 100 |./samh 18| wc
```

2) Ancient population expansion:

```
./ms 178 100000 -s 220 -r 11610.4 29027 -c 1 100 -G 55262.04223 -eG 0.000125 0 |./samh 18| wc
```

3) Recent population expansion:

```
./ms 178 100000 -s 220 -r 11610.4 29027 -c 1 100 -G 276310.2112 -eG 0.000025 0 |./samh 18| wc
```

4) Several bottleneck:

```
./ms 178 100000 -s 220 -r 11610.4 29027 -c 1 100 -G 147365.446 -eG 0.0000625 0 -eN 0.000125 0.001 |./samh 18| wc
```

5) Repeated bottlenecks with subsequent expansion:

```
./ms 178 100000 -s 220 -r 11610.4 29027 -c 1 100 -G 147365.446 -eG 0.0000625 0 -eN 0.000075 0.001 -eG 0.000075 184206.8074 -eG 0.0000875 0 -eN 0.0001 0.001 -eG 0.0001 184206.8074 -eG 0.0001125 0 -eN 0.000125 0.001 -eG 0.000125 184206.8074 -eG 0.0001375 0 -eN 0.00015 0.001 -eG 0.00015 184206.8074 -eG 0.0001625 0 -eN 0.000175 0.001 |./samh 18| wc
```

6) Population structure with 2 subpopulations:

```
./ms 178 100000 -s 220 -r 11.6104 29027 -c 1 100 -es 0.0 1 .5 -eM 0.0 1.0 |./samh 18| wc
```

7) Population structure with 3 subpopulations:

```
./ms 178 100000 -s 220 -r 11.6104 29027 -c 1 100 -es 0.0 1 0.3333 -es 0.0 1 0.5 -eM 0.0 1.0 |./samh 18| wc
```

8) Population structure with 4 subpopulations:

```
./ms 178 100000 -s 220 -r 11.6104 29027 -c 1 100 -es 0.0 1 .25 -es 0.0 1 .333 -es 0.0 1 .5 -eM 0.0 1.0 |./samh 18| wc
```

9) Population structure with 5 subpopulations:

```
./ms 178 100000 -s 220 -r 11.6104 29027 -c 1 100 -es 0.0 1 0.2 -es 0.0 1 0.25 -es 0.0 1 0.333 -es 0.0 1 0.5 -eM 0.0 1.0 |./samh 18| wc
```

### *Age estimation*

We estimated the age of haplogroup D using a mutation-based method as previously described (S11). This method simply relies on averaging the number of mutations along each lineage from the most recent common ancestor (MRCA) to the sampled chromosome. This averaging produces an estimate of the time to MRCA that is unbiased by demographic history (S11). Let  $t$  denote the time to MRCA for haplogroup D in units of mutations. The value of  $t$  could be estimated as follows: To start, we decided to focus only on the 47 individuals who are homozygous for haplogroup D chromosomes (rather than using all the inferred copies of haplogroup D). This avoided uncertainties in haplotype inference. We also note that there are no evident recombinants between D and non-D types among these 47 individuals, which is important because the absence of such recombinants is a necessary condition for our methodology (S11, S12). Using chimpanzee sequence as an outgroup, we deduced the MRCA sequence of haplogroup D, which happens to be the same as the sequence of haplotype 49. We next added up the total number of mutations separating the MRCA and the 94 chromosomes sampled in the 47 individuals. This number was 93, which was divided by 94 to yield  $\hat{t}$ , the estimate of  $t$ , at 0.989. This value was then divided by 23416 (the total length of DNA sequenced) to yield an estimate for the number of mutations per base ( $\hat{T}$ ) of  $4.2 \times 10^{-5}$ . By comparing human and chimpanzee sequences in this region, the rate of human-chimpanzee nucleotide divergence ( $D$ ) in this region was estimated at 0.0136 mutations per base. Finally, human-chimpanzee divergence time ( $L$ ) was set at  $6 \times 10^6$  years. Most estimates of this time is between  $5 \times 10^6$

and  $6 \times 10^6$  years. We chose the upper one to be conservative. The estimated time to MRCA in years was then obtained, using the simple formula  $(2\hat{T}/D)*L$  as described previously (S11), at 37,281 years before present. The coalescence age of the entire Coriell panel was calculated in a similar manner. There are a total of 8136 mutations between the 178 chromosomes in the Coriell panel and the deduced MRCA sequence, which leads to an age estimate of 1,722,347 years. We note that owing to recombination, this estimated age is actually the average of multiple coalescence ages corresponding to multiple recombination blocks that coalesce independently.

The 95% confidence interval (CI) for the age of haplogroup D was estimated by an analytical approach that is an extension of a previously described method (S12). Let  $y_i$  denote the number of differences between the MRCA and the  $i^{\text{th}}$  chromosome. The value of  $\hat{t}$  would be  $(\sum_{i=1}^n y_i)/n$ , where  $n$  is the number of chromosomes sampled. The variance of  $\hat{t}$  is  $[\sum_{i=1}^n \text{var}(y_i) + 2\sum_{i<j} \text{cov}(y_i, y_j)]/n^2$ . If we assume an infinite-sites model, each  $y_i$  is Poisson distributed with mean  $t$ . The  $\text{var}(y_i)$  is simply  $t$ , and the  $\text{cov}(y_i, y_j)$  is simply  $t - t_{ij}$ , where  $t_{ij}$  is the time of the most recent common ancestor of chromosome  $i$  and chromosome  $j$  (S11, S12). Therefore the variance of our estimate is  $t/n + 2[\sum_{i<j} (t - t_{ij})]/n^2$ . There are  $n(n-1)/2$  terms in this sum, so this can be written as  $t - 2[\sum_{i<j} (t_{ij})]/n^2$  or  $t - [(n-1)/n]\bar{t}_{ij}$ , where  $\bar{t}_{ij}$  is the average time to the most recent common ancestor of a pair of chromosomes.  $\bar{t}_{ij}$  can be estimated as one-half the average pairwise differences between the 94 chromosomes, calculated as  $(1/2)\sum_{k=1}^m \{2f_k(m - f_k)/[m(m-1)]\}$  or  $\sum_{k=1}^m \{f_k(m - f_k)/[m(m-1)]\}$ , where  $f_k$  is the count of the derived allele at the  $k^{\text{th}}$  polymorphic site and  $m$  is the total number of polymorphic sites. So we can estimate the variance of  $\hat{t}$  by  $\hat{t} - [(n-1)/n]\sum_{k=1}^m \{f_k(m - f_k)/[m(m-1)]\}$ . For the 94 haplogroup D chromosomes sampled in the 47 individuals, there are 34 SNP sites. Let  $N_x$  designate the number of sites where the count of the derived allele is  $x$ . For our data,  $N_1 = 23$ ,  $N_2 = 2$ ,  $N_3 = 5$ ,  $N_4 = 1$ ,  $N_{15} = 2$ ,  $N_{17} = 1$ , and all others  $N_x$  values are zero. Thus, based on our data, the estimate for the variance of  $\hat{t}$  is 0.094, and the estimate for the standard error of  $\hat{t}$  is  $\sqrt{0.094} = 0.307$ . Assuming that the  $\hat{t}$  estimator is roughly normally distributed, the 95% CI of  $\hat{t}$  would be approximately 0.376 to 1.60. This corresponds, in units of years, a CI of 14175 to 60387 years before present. We note that this CI does not consider uncertainties in mutation rate. It also does not consider uncertainties in the estimated human-chimpanzee divergence time, which can only be inferred from fossil records and molecular data, and cannot be directly observed.



**References and Notes**

- S1. M. Stephens, N. J. Smith, P. Donnelly, *Am. J. Hum. Genet.* **68**, 978 (2001).
- S2. M. Stephens, P. Donnelly, *Am. J. Hum. Genet.* **73**, 1162 (2003).
- S3. J. Rozas, R. Rozas, *Bioinformatics* **15**, 174 (1999).
- S4. H. M. Cann *et al.*, *Science* **296**, 261 (2002).
- S5. B. S. Weir, C. C. Cockerham, *Evolution* **38**, 1358 (1984).
- S6. B. S. Weir, *Genetic Data Analysis* (Sinauer Associates, Sunderland, 1990).
- S7. R. R. Hudson, *Oxford Surv.Evol. Biol.* **7**, 1 (1990).
- S8. R. R. Hudson, *Bioinformatics* **18**, 337 (2002).
- S9. A. Kong *et al.*, *Nat. Genet.* **31**, 241 (2002).
- S10. C. Wiuf, J. Hein, *Genetics* **155**, 451 (2000).
- S11. R. Thomson, J. K. Pritchard, P. Shen, P. J. Oefner, M. W. Feldman, *Proc. Natl. Acad. Sci. U.S.A.* **97**, 7360 (2000).
- S12. H. Tang, D. O. Siegmund, P. Shen, P. J. Oefner, M. W. Feldman, *Genetics* **161**, 447 (2002).











Supplementary Table S2. Haplotype data of *Microcephalin* in the 89 individuals of the Coriell panel.

Haplotype	Occurrence	Genomic position from start codon																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																																													
		29022	29134	29250	29366	29482	29598	29714	29830	29946	30062	30178	30294	30410	30526	30642	30758	30874	30990	31106	31222	31338	31454	31570	31686	31802	31918	32034	32150	32266	32382	32498	32614	32730	32846	32962	33078	33194	33310	33426	33542	33658	33774	33890	34006	34122	34238	34354	34470	34586	34702	34818	34934	35050	35166	35282	35398	35514	35630	35746	35862	35978	36094	36210	36326	36442	36558	36674	36790	36906	37022	37138	37254	37370	37486	37602	37718	37834	37950	38066	38182	38298	38414	38530	38646	38762	38878	38994	39110	39226	39342	39458	39574	39690	39806	39922	40038	40154	40270	40386	40502	40618	40734	40850	40966	41082	41198	41314	41430	41546	41662	41778	41894	42010	42126	42242	42358	42474	42590	42706	42822	42938	43054	43170	43286	43402	43518	43634	43750	43866	43982	44098	44214	44330	44446	44562	44678	44794	44910	45026	45142	45258	45374	45490	45606	45722	45838	45954	46070	46186	46302	46418	46534	46650	46766	46882	46998	47114	47230	47346	47462	47578	47694	47810	47926	48042	48158	48274	48390	48506	48622	48738	48854	48970	49086	49202	49318	49434	49550	49666	49782	49898	50014	50130	50246	50362	50478	50594	50710	50826	50942	51058	51174	51290	51406	51522	51638	51754	51870	51986	52102	52218	52334	52450	52566	52682	52798	52914	53030	53146	53262	53378	53494	53610	53726	53842	53958	54074	54190	54306	54422	54538	54654	54770	54886	55002	55118	55234	55350	55466	55582	55698	55814	55930	56046	56162	56278	56394	56510	56626	56742	56858	56974	57090	57206	57322	57438	57554	57670	57786	57902	58018	58134	58250	58366	58482	58598	58714	58830	58946	59062	59178	59294	59410	59526	59642	59758	59874	59990	60106	60222	60338	60454	60570	60686	60802	60918	61034	61150	61266	61382	61498	61614	61730	61846	61962	62078	62194	62310	62426	62542	62658	62774	62890	63006	63122	63238	63354	63470	63586	63702	63818	63934	64050	64166	64282	64398	64514	64630	64746	64862	64978	65094	65210	65326	65442	65558	65674	65790	65906	66022	66138	66254	66370	66486	66602	66718	66834	66950	67066	67182	67298	67414	67530	67646	67762	67878	67994	68110	68226	68342	68458	68574	68690	68806	68922	69038	69154	69270	69386	69502	69618	69734	69850	69966	70082	70198	70314	70430	70546	70662	70778	70894	71010	71126	71242	71358	71474	71590	71706	71822	71938	72054	72170	72286	72402	72518	72634	72750	72866	72982	73098	73214	73330	73446	73562	73678	73794	73910	74026	74142	74258	74374	74490	74606	74722	74838	74954	75070	75186	75302	75418	75534	75650	75766	75882	75998	76114	76230	76346	76462	76578	76694	76810	76926	77042	77158	77274	77390	77506	77622	77738	77854	77970	78086	78202	78318	78434	78550	78666	78782	78898	79014	79130	79246	79362	79478	79594	79710	79826	79942	80058	80174	80290	80406	80522	80638	80754	80870	80986	81102	81218	81334	81450	81566	81682	81798	81914	82030	82146	82262	82378	82494	82610	82726	82842	82958	83074	83190	83306	83422	83538	83654	83770	83886	84002	84118	84234	84350	84466	84582	84698	84814	84930	85046	85162	85278	85394	85510	85626	85742	85858	85974	86090	86206	86322	86438	86554	86670	86786	86902	87018	87134	87250	87366	87482	87598	87714	87830	87946	88062	88178	88294	88410	88526	88642	88758	88874	88990	89106	89222	89338	89454	89570	89686	89802	89918	90034	90150	90266	90382	90498	90614	90730	90846	90962	91078	91194	91310	91426	91542	91658	91774	91890	92006	92122	92238	92354	92470	92586	92702	92818	92934	93050	93166	93282	93398	93514	93630	93746	93862	93978	94094	94210	94326	94442	94558	94674	94790	94906	95022	95138	95254	95370	95486	95602	95718	95834	95950	96066	96182	96298	96414	96530	96646	96762	96878	96994	97110	97226	97342	97458	97574	97690	97806	97922	98038	98154	98270	98386	98502	98618	98734	98850	98966	99082	99198	99314	99430	99546	99662	99778	99894	100010	100126	100242	100358	100474	100590	100706	100822	100938	101054	101170	101286	101402	101518	101634	101750	101866	101982	102098	102214	102330	102446	102562	102678	102794	102910	103026	103142	103258	103374	103490	103606	103722	103838	103954	104070	104186	104302	104418	104534	104650	104766	104882	104998	105114	105230	105346	105462	105578	105694	105810	105926	106042	106158	106274	106390	106506	106622	106738	106854	106970	107086	107202	107318	107434	107550	107666	107782	107898	108014	108130	108246	108362	108478	108594	108710	108826	108942	109058	109174	109290	109406	109522	109638	109754	109870	109986	110102	110218	110334	110450	110566	110682	110798	110914	111030	111146	111262	111378	111494	111610	111726	111842	111958	112074	112190	112306	112422	112538	112654	112770	112886	113002	113118	113234	113350	113466	113582	113698	113814	113930	114046	114162	114278	114394	114510	114626	114742	114858	114974	115090	115206	115322	115438	115554	115670	115786	115902	116018	116134	116250	116366	116482	116598	116714	116830	116946	117062	117178	117294	117410	117526	117642	117758	117874	117990	118106	118222	118338	118454	118570	118686	118802	118918	119034	119150	119266	119382	119498	119614	119730	119846	119962	120078	120194	120310	120426	120542	120658	120774	120890	121006	121122	121238	121354	121470	121586	121702	121818	121934	122050	122166	122282	122398	122514	122630	122746	122862	122978	123094	123210	123326	123442	123558	123674	123790	123906	124022	124138	124254	124370	124486	124602	124718	124834	124950	125066	125182	125298	125414	125530	125646	125762	125878	125994	126110	126226	126342	126458	126574	126690	126806	126922	127038	127154	127270	127386	127502	127618	127734	127850	127966	128082	128198	128314	128430	128546	128662	128778	128894	129010	129126	129242	129358	129474	129590	129706	129822	129938	130054	130170	130286	130402	130518	130634	130750	130866	130982	131098	131214	131330	131446	131562	131678	131794	131910	132026	132142	132258	132374	132490	132606	132722	132838	132954	133070	133186	133302	133418	133534	133650	133766	133882	133998	134114	134230	134346	134462	134578	134694	134810	134926	135042	135158	135274	135390	135506	135622	135738	135854	135970	136086	136202	136318	136434	136550	136666	136782	136898	137014	137130	137246	137362	137478	137594	137710	137826	137942	138058	138174	138290	138406	138522	138638	138754	138870	138986	139102	139218	139334	139450	139566	139682	139798	139914	140030	140146	140262	140378	140494	140610	140726	140842	140958	141074	141190	141306	141422	141538	141654	141770	141886	142002	142118	142234	142350	142466	142582	142698	142814	142930	143046	143162	143278	143394	143510	143626	143742	143858	143974	144090	144206	144322	144438	144554	144670	144786	144902	145018	145134	145250	145366	145482	145598	145714	145830	145946	146062	146178	146294	146410	146526	146642	146758	146874	146990	147106	147222	147338	147454	147570	147686	147802	147918	148034	148150	148266	148382	148498	148614	148730	148846	148962	149078	149194	149310	149426	149542	149658	149774	149890	150006	150122	150238	150354	150470	150586	150702	150818	150934	151050	151166	151282	151398	151514	151630	151746	151862	151978	152094	152210	152326	152442	152558	152674	152790	152906	153022	153138	153254	153370	153486	153602	153718	153834	153950	154066	154182	154298	154414	154530	154646	154762	154878	154994	155110	155226	155342	155458	155574	155690	155806	155922	156038	156154	156270	156386	156502	156618	156734	156850	156966	157082	157198	157314	157430	157546	157662	157778	157894	158010	158126	158242	158358	158474	158590	158706	158822	158938	159054	159170	159286	159402	159518	159634	159750	159866	159982	160098	160214	160330	160446	160562	160678	160794	160910	161026	161142	161258	161374	161490	161606	161722	161838	161954	162070	162186	162302	162418	162534	162650	162766	162882	162998	163114	163230	163346	163462	163578	163694	163810	163926	164042	164158





