

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 (π) 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 π and Tajima's D of this haplogroup. Inferred haplotypes were used to calculate π 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 4.2×10^{-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 6×10^6 years. Most estimates of this time is between 5×10^6

and 6×10^6 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^{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^{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).

