

# Supporting Information

Xu et al. 10.1073/pnas.1118892109

## SI Materials and Methods

**Populations and Samples.** The Pan-Asian dataset consists of one Papuan population from Papua New Guinea and one Melanesian population from Bougainville obtained from the database of the Human Genome Diversity Project–Centre d'Etude du Polymorphisme Humain and 288 unrelated individuals representing 13 Indonesian populations obtained from the Pan-Asia SNP Project (1) (Table 1). The Affymetrix 6.0 dataset consists of 61 individuals from seven populations from Indonesia and one Papuan population from the southern highlands of Papua New Guinea (2, 3). Detailed information on these populations is presented in Table 2.

**Markers and Their Positions.** In the Pan-Asian dataset, genotype data of 13 Indonesian populations generated using Affymetrix Genechip Human Mapping 50K Xba array were obtained from the Pan-Asia SNP Project (1). Detailed information about data generation, filtration, and data quality control was described elsewhere (1). Genotype data of one Papuan population and one Melanesian population were generated using Illumina Genechip Human Mapping 650K Xba array, and details were described elsewhere (4). All of the analyses in this study used the markers that genotyped in both Pan-Asia and HGDP-CEPH samples. With data integration, we obtained 19,934 SNPs shared by 15 population samples. The physical positions of SNPs were based on the *Homo sapiens* Genome Build 37. The average spacing between adjacent markers was 50.5 kb, with a minimum of 11 bp and a maximum of 26.8 Mb; the median between marker distances was 18.3 kb. In the Affymetrix 6.0 dataset, all samples were genotyped on the Affymetrix 6.0 platform as described previously (2, 3). After data cleaning and integration, there were 685,582 SNPs for analysis.

**Statistical Analysis. Analysis of the Pan-Asian dataset.** Principle component analysis was performed at the individual level using EIGENSOFT version 3.0 (5). Unbiased estimates of  $F_{ST}$  were calculated using the work by Weir and Hill (6) with PEAS V1.0 (7). Great circle distance calculations followed the approach in the work by Ramachandran et al. (8). The tree of populations was reconstructed based on the  $F_{ST}$  distances and the neighbor-joining algorithm (9) implemented in the Molecular Evolutionary Genetics Analysis software package (MEGA version 4.0) (10).

Given the large number of markers in our dataset, genetic analyses can be performed at the level of individual, making no presumption of group membership. We applied a Bayesian cluster analysis as implemented in the STRUCTURE program (11) and a maximum likelihood method as implemented in the *frappe* program (12) to infer the genetic ancestry of individuals. Our approach is solely based on genotype without incorporating any

information on sampling location or population affiliation of each individual. We ran STRUCTURE from  $K = 1$  to  $K = 15$ , and 10 repeats were done for each  $K$  values. All STRUCTURE runs used 10,000 iterations after a burn-in of length 20,000 with the admixture model and assuming that allele frequencies were correlated (11). The *frappe* program was run for 100,000 iterations from  $K = 2$ –15 and repeated 10 times for each single  $K$  value. According to the distribution of the posterior probability as provided by STRUCTURE and *frappe* analyses, the most probable and appropriate number of clusters should be three in our dataset.

In estimating the admixture time of East Indonesian populations, we selected a panel of 2,807 ancestry informative markers with large allele frequency differences ( $F_{ST} > 0.3$ ) between ID-MT and Papuan and ran STRUCTURE with the linkage model to estimate recombination rates in seven Eastern Indonesian populations. In this model, STRUCTURE reports not only the overall ancestry for each individual but also the probability of origin of each allele. The break points were inferred according to the estimated origin of each allele. The program STRUCTURE was run with 100,000 iterations, 200,000 burn-ins, and 10,000 admixture burn-ins.

**Analysis of the Affymetrix 6.0 dataset.** Principle component analysis, admixture proportions, and time of admixture estimation were performed using the StepPCO software (13). Individual ancestry components were inferred using a maximum likelihood method as implemented in the *frappe* program (12). We ran analyses for  $K = 2$  and  $K = 3$  and performed three independent runs for each  $K$  value. The analysis of admixture rates on the autosomes vs. X chromosome was based on 36,415 X-linked SNPs. For the time of admixture estimation, Borneo and New Guinea were used as parental groups. To make the sample sizes equal for the two parental groups, 16 individuals were selected at random from the NGH population. To limit the analysis only to variation defined by the parental groups and exclude any signal in the admixed groups that comes from genetic drift or other sources of admixture, the first principle axis was calculated only between the parental groups, and the admixed group was then projected onto this axis (14, 15). The admixture signal along each chromosome was obtained, and the width of the ancestry blocks was estimated as described previously (13). The method is sensitive to small sample sizes, and the admixed Indonesian groups were, therefore, combined into the Nusa Tenggara group (10 individuals from the islands of Alor, Timor, Roti, and Flores) and the Moluccas (10 individuals from the islands of Hiri and Ternate). The time estimate is based on comparison with the data from 100 forward simulations with a 40% migration rate (13).

1. Abdulla MA, et al. (2009) Mapping human genetic diversity in Asia. *Science* 326: 1541–1545.
2. Reich D, et al. (2011) Denisova admixture and the first modern human dispersals into Southeast Asia and Oceania. *Am J Hum Genet* 89:516–528.
3. Wollstein A, et al. (2010) Demographic history of Oceania inferred from genome-wide data. *Curr Biol* 20:1983–1992.
4. Li JZ, et al. (2008) Worldwide human relationships inferred from genome-wide patterns of variation. *Science* 319:1100–1104.
5. Patterson N, Price AL, Reich D (2006) Population structure and eigenanalysis. *PLoS Genet* 2:e190.
6. Weir BS, Hill WG (2002) Estimating F-statistics. *Annu Rev Genet* 36:721–750.
7. Xu S, Gupta S, Jin L (2010) PEAS V1.0: A package for elementary analysis of SNP data. *Mol Ecol Resour* 10:1085–1088.
8. Ramachandran S, et al. (2005) Support from the relationship of genetic and geographic distance in human populations for a serial founder effect originating in Africa. *Proc Natl Acad Sci USA* 102:15942–15947.
9. Saitou N, Nei M (1987) The neighbor-joining method: A new method for reconstructing phylogenetic trees. *Mol Biol Evol* 4:406–425.
10. Tamura K, Dudley J, Nei M, Kumar S (2007) MEGA4: Molecular Evolutionary Genetics Analysis (MEGA) software version 4.0. *Mol Biol Evol* 24:1596–1599.
11. Falush D, Stephens M, Pritchard JK (2003) Inference of population structure using multilocus genotype data: Linked loci and correlated allele frequencies. *Genetics* 164: 1567–1587.
12. Tang H, Peng J, Wang P, Risch NJ (2005) Estimation of individual admixture: Analytical and study design considerations. *Genet Epidemiol* 28:289–301.
13. Pugach I, Matveyev R, Wollstein A, Kayser M, Stoneking M (2011) Dating the age of admixture via wavelet transform analysis of genome-wide data. *Genome Biol* 12:R19.
14. Reich D, Thangaraj K, Patterson N, Price AL, Singh L (2009) Reconstructing Indian population history. *Nature* 461:489–494.
15. McVean G (2009) A genealogical interpretation of principal components analysis. *PLoS Genet* 5:e1000686.

## SI Text

The participants of the HUGO Pan-Asian SNP Consortium are arranged by surname alphabetically in the following:

Mahmood Ameen Abdulla<sup>a</sup>, Ikhlaq Ahmed<sup>b</sup>, Anunchai Assawamakin<sup>c,d</sup>, Jong Bhak<sup>e</sup>, Samir K. Brahmachari<sup>b</sup>, Gayvelline C. Calacal<sup>f</sup>, Amit Chaurasia<sup>b</sup>, Chien-Hsiun Chen<sup>g</sup>, Jieming Chen<sup>h</sup>, Yuan-Tsong Chen<sup>g</sup>, Jiayou Chu<sup>i</sup>, Eva Maria C. Cutiongco-de la Paz<sup>j</sup>, Maria Corazon A. De Ungria<sup>f</sup>, Frederick C. Delfin<sup>f</sup>, Juli Edo<sup>a</sup>, Suthat Fuchareon<sup>c</sup>, Ho Ghang<sup>e</sup>, Takashi Gojobori<sup>k,l</sup>, Junsong Han<sup>m</sup>, Sheng-Feng Ho<sup>g</sup>, Boon Peng Hoh<sup>n</sup>, Wei Huang<sup>o</sup>, Hidetoshi Inoko<sup>p</sup>, Pankaj Jha<sup>b</sup>, Timothy A. Jinam<sup>1</sup>, Li Jin<sup>q,r</sup>, Jongsun Jung<sup>s</sup>, Daoroong Kangwanpong<sup>t</sup>, Jatupol Kampuansai<sup>t</sup>, Giulia C. Kennedy<sup>u,v</sup>, Preeti Khurana<sup>w</sup>, Hyung-Lae Kim<sup>s</sup>, Kwangjoong Kim<sup>s</sup>, Sangsoo Kim<sup>x</sup>, Woo-Yeon Kim<sup>e</sup>, Kuchan Kimm<sup>y</sup>, Ryosuke Kimura<sup>z</sup>, Tomohiro Koike<sup>k</sup>, Supasak Kulawonganunчай<sup>d</sup>, Vikrant Kumar<sup>h</sup>, Poh San Lai<sup>aa,bb</sup>, Jong-Young Lee<sup>s</sup>, Sunghoon Lee<sup>e</sup>, Edison T. Liu<sup>h</sup>, Partha P. Majumder<sup>cc</sup>, Kiran Kumar Mandapati<sup>ww</sup>, Sangkot Marzuki<sup>dd</sup>, Wayne Mitchell<sup>ee,ff</sup>, Mitali Mukerji<sup>b</sup>, Kenji Naritomi<sup>gg</sup>, Chumpol Ngamphiw<sup>d</sup>, Norio Niikawa<sup>hh</sup>, Nao Nishida<sup>z</sup>, Bermseok Oh<sup>s</sup>, Sangho Oh<sup>e</sup>, Jun Ohashi<sup>z</sup>, Akira Oka<sup>p</sup>, Rick Ong<sup>h</sup>, Carmencita D. Padilla<sup>j</sup>, Prasit Palittapongarnpim<sup>ii</sup>, Henry B. Perdigon<sup>f</sup>, Maude Elvira Phipps<sup>ajj</sup>, Eileen Png<sup>h</sup>, Yoshiyuki Sakaki<sup>kk</sup>, Jazelyn M. Salvador<sup>f</sup>, Yuliana Sandraling<sup>dd</sup>, Vinod Scaria<sup>b</sup>, Mark Seielstad<sup>h</sup>, Mohd Ros Sidek<sup>n</sup>, Amit Sinha<sup>b</sup>, Metawee Srikumool<sup>t</sup>, Herawati Sudoyo<sup>dd</sup>, Sumio Sugano<sup>ll</sup>, Helena Suryadi<sup>dd</sup>, Yoshiyuki Suzuki<sup>k</sup>, Kristina A. Tabbada<sup>f</sup>, Adrian Tan<sup>h</sup>, Katsushi Tokunaga<sup>z</sup>, Sissades Tongsimad<sup>d</sup>, Lilian P. Villamor<sup>f</sup>, Eric Wang<sup>uu,v</sup>, Ying Wang<sup>o</sup>, Haifeng Wang<sup>o</sup>, Jer-Yuarn Wu<sup>g</sup>, Huasheng Xiao<sup>mm</sup>, Shuhua Xu<sup>r</sup>, Jin Ok Yang<sup>e</sup>, Yin Yao Shugart<sup>mmm</sup>, Hyang-Sook Yoo<sup>e</sup>, Wentao Yuan<sup>o</sup>, Guoping Zhao<sup>o</sup>, Bin Alwi Zilfalil<sup>n</sup>, and Indian Genome Variation Consortium<sup>b</sup>

<sup>a</sup>Department of Molecular Medicine, Faculty of Medicine, and Department of Anthropology, Faculty of Arts and Social Sciences, University of Malaya, Kuala Lumpur, 50603, Malaysia; <sup>b</sup>Council for Scientific and Industrial Research, Institute of Genomics and Integrative Biology, Delhi 110007, India; <sup>c</sup>Mahidol University, Salaya Campus, Puttamonthon, Nakornpathom 73170, Thailand; <sup>d</sup>Biostatistics and Informatics Laboratory, Genome Institute, National Center for Genetic Engineering and Biotechnology,

Pathumtani 12120, Thailand; <sup>e</sup>Korean BioInformation Center (KOBIC), Korea Research Institute of Bioscience and Biotechnology (KRIBB), Yuseong-gu, Daejeon 305-806, Korea; <sup>f</sup>DNA Analysis Laboratory, Natural Sciences Research Institute, University of the Philippines, Diliman, Quezon City 1101, Philippines; <sup>g</sup>Institute of Biomedical Sciences, Academia Sinica, Nangang, Taipei City 115, Taiwan; <sup>h</sup>Genome Institute of Singapore, 138672, Singapore; <sup>i</sup>Institute of Medical Biology, Chinese Academy of Medical Science, Kunming 650118, China; <sup>j</sup>Institute of Human Genetics, National Institutes of Health, University of the Philippines Manila, Ermita Manila 1000, Philippines; <sup>k</sup>Center for Information Biology and DNA Data Bank of Japan, National Institute of Genetics, Research Organization of Information and Systems, Mishima, Shizuoka 411-8540, Japan; <sup>l</sup>Biomedical Information Research Center, National Institute of Advanced Industrial Science and Technology, Koto-ku, Tokyo 135-0064, Japan; <sup>m</sup>National Engineering Center for Biochip at Shanghai, Shanghai 201203, China; <sup>n</sup>Human Genome Center, School of Medical Sciences, Universiti Sains Malaysia, 16150 Kubang Kerian, Kelantan, Malaysia; <sup>o</sup>Ministry of Science and Technology-Shanghai Laboratory of Disease and Health Genomics, Chinese National Human Genome Center Shanghai, Shanghai 201203, China; <sup>p</sup>Division of Molecular Medical Science and Molecular Medicine, Department of Molecular Life Science, Tokai University School of Medicine, Isehara-A Kanagawa-Pref A259-1193, Japan; <sup>q</sup>State Key Laboratory of Genetic Engineering and Ministry of Education Key Laboratory of Contemporary Anthropology, School of Life Sciences, Fudan University, Shanghai 200433, China; <sup>r</sup>Chinese Academy of Sciences-Max Planck Society Partner Institute for Computational Biology, Shanghai Institutes of Biological Sciences, Chinese Academy of Sciences, Shanghai 200031, China; <sup>s</sup>Korea National Institute of Health, Eunpyung-Gu, Seoul, 122-701, Korea; <sup>t</sup>Department of Biology, Faculty of Science, Chiang Mai University, Chiang Mai 50202, Thailand; <sup>u</sup>Genomics Collaborations, Affymetrix, Santa Clara, CA 95051; <sup>v</sup>Veracyte, South San Francisco, CA 94080; <sup>w</sup>The Centre for Genomic Applications (an Institute of Genomics and Integrative Biology-Institute of Molecular Medicine Collaboration), New Delhi 110020, India; <sup>x</sup>Soongsil University, Dongjak-gu, Seoul 156-743, Korea; <sup>y</sup>Eulji University College of Medicine, Dae-jeon City 301-832, Korea; <sup>z</sup>Department of Human Genetics, Graduate School of Medicine, University of Tokyo, Bunkyo-ku, Tokyo 113-0033, Japan; <sup>aa</sup>Department of Pediatrics, Yong Loo Lin School of Medicine, National University of Singapore, National University Hospital, 119074, Singapore; <sup>bb</sup>Population Genetics Laboratory, Defense Medical and Environmental Research Institute, Defence Science Organisation National Laboratories, 117510, Singapore; <sup>cc</sup>Indian Statistical Institute (Kolkata), Kolkata 700108, India; <sup>dd</sup>Eijkman Institute for Molecular Biology, Jakarta 10430, Indonesia; <sup>ee</sup>Informatics Experimental Therapeutic Centre, 03-01 Nanos, 138669, Singapore; <sup>ff</sup>Division of Information Sciences, School of Computer Engineering, Nanyang Technological University, 639798, Singapore; <sup>gg</sup>Department of Medical Genetics, Faculty of Medicine, University of the Ryukyus, Nishihara, Okinawa 903-0215, Japan; <sup>hh</sup>Research Institute of Personalized Health Sciences, Health Sciences University of Hokkaido, Tobetsu 061-0293, Japan; <sup>ii</sup>National Science and Technology Development Agency, Pathumtani 12120, Thailand; <sup>jj</sup>Monash University (Sunway Campus), 46150 Bandar Sunway, Selangor, Malaysia; <sup>kk</sup>RIKEN Genomic Sciences Center, Tsurumi-ku, Yokohama 230-0045, Japan; <sup>ll</sup>Laboratory of Functional Genomics, Department of Medical Genome Sciences Graduate School of Frontier Sciences, Shirokanedai Laboratory, University of Tokyo, Minato-ku, Tokyo 108-8639, Japan; and <sup>mmm</sup>Genomic Research Branch, National Institute of Mental Health, National Institutes of Health, Bethesda, MD 20892

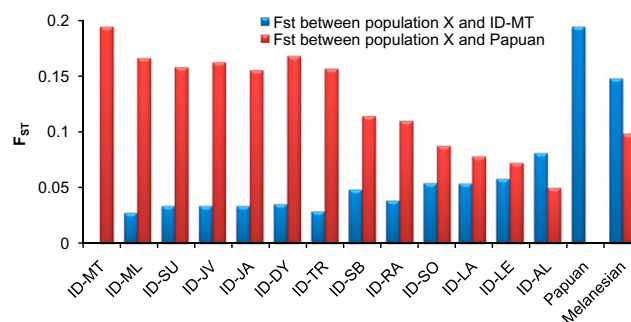


Fig. S1.  $F_{ST}$  clines for 15 population samples. Pairwise  $F_{ST}$  values of all populations in the Pan-Asian SNP dataset were compared with the easternmost (Papuan) and westernmost (ID-MT) populations, respectively.





**Table S2. Admixture proportion of populations (%) estimated from 19,934 SNPs**

	Cluster 1	Cluster 2	Cluster 3
ID-MT	99.9	0.0	0.0
ID-ML	98.5	1.0	0.5
ID-SU	98.9	0.7	0.4
ID-JV	99.2	0.6	0.2
ID-JA	99.4	0.4	0.2
ID-DY	99.2	0.5	0.4
ID-TR	94.7	5.1	0.2
ID-SB	76.2	23.7	0.1
ID-RA	75.6	24.1	0.3
ID-SO	64.9	34.9	0.3
ID-LA	60.7	39.0	0.2
ID-LE	57.7	41.9	0.5
ID-AL	43.8	55.4	0.9
Papuan	1.2	98.4	0.4
Melanesian	6.7	47.3	46.0

Cluster 1, Asian; Cluster 2, Papuan; Cluster 3, Melanesian. Note that admixture proportions in the table are the results averaged from 10 independent structure runs; the variation of the estimations from different runs is very small, and SDs are less than 1%.