

# 1 **Multi-layered population structure in Island Southeast Asians**

2  
3 **Running title: Population structure in Island Southeast Asia**

4  
5 Alexander Mörseburg<sup>1+\*</sup>, Luca Pagani<sup>1,2\*</sup>, Francois-Xavier Ricaut<sup>3</sup>, Bryndis Yngvadottir<sup>1</sup>,  
6 Eadaoin Harney<sup>1</sup>, Cristina Castillo<sup>4</sup>, Tom Hoogervorst<sup>5</sup>, Tiago Antao<sup>6</sup>, Pradiptajati Kusuma<sup>3,7</sup>,  
7 Nicolas Brucato<sup>3</sup> Alexia Cardona<sup>1</sup>, Denis Pierron<sup>3</sup>; Thierry Letellier<sup>3</sup>, Joseph Wee<sup>8</sup>, Syafiq  
8 Abdullah<sup>9</sup>, Mait Metspalu<sup>10,11</sup>, Toomas Kivisild<sup>1,10</sup>

9  
10 <sup>1</sup>Division of Biological Anthropology, University of Cambridge, Fitzwilliam Street,  
11 Cambridge CB2 1QH, United Kingdom.

12 <sup>2</sup>Department of Biological, Geological and Environmental Sciences, University of Bologna,  
13 Via Selmi 3, 40126 Bologna, Italy.

14 <sup>3</sup>Laboratoire d'Anthropologie Moléculaire et Imagerie de Synthèse, UMR 5288, Centre  
15 National de la Recherche Scientifique, Université de Toulouse, 31073 Toulouse, France

16 <sup>4</sup>Institute of Archaeology, University College London, 31-34 Gordon Square, London WC1H  
17 0PY, UK

18 <sup>5</sup>Royal Netherlands Institute of Southeast Asian and Caribbean Studies, Reuvenplaats 2,  
19 2311 BE Leiden, Netherlands

20 <sup>6</sup>Department of Vector Biology, Liverpool School of Tropical Medicine, Pembroke Place  
21 Liverpool L3 5QA, United Kingdom

22 <sup>7</sup>Genome Diversity and Diseases Laboratory, Eijkman Institute for Molecular Biology, Jl.  
23 Diponegoro 69, Jakarta 10430, Indonesia.

24 <sup>8</sup>National Cancer Centre, Singapore 169610, Singapore;

25 <sup>9</sup>RIPAS Hospital, Bandar Seri Begawan, BA 1710, Brunei Darussalam

26 <sup>10</sup>Evolutionary Biology Group, Estonian Biocentre, 51010 Tartu, Estonia;

27 <sup>11</sup>Department of Evolutionary Biology, Institute of Molecular and Cell Biology, University of  
28 Tartu, 51010 Tartu, Estonia

29  
30 \* AM and LP contributed equally to the work

31 <sup>+</sup>Corresponding author, E-mail address: [am2037@cam.ac.uk](mailto:am2037@cam.ac.uk), +44 7900 373 200

32  
33  
34  
35

1 **ABSTRACT**

2 The history of human settlement in Southeast Asia has been complex and involved several  
3 distinct dispersal events. Here we report the analyses of 1825 individuals from Southeast Asia  
4 including new genome-wide genotype data for 146 individuals from three Mainland Southeast  
5 Asian (Burmese, Malay and Vietnamese) and four Island Southeast Asian (Dusun, Filipino,  
6 Kankanaey and Murut) populations. While confirming the presence of previously recognized  
7 major ancestry components in the Southeast Asian population structure, we highlight the  
8 Kankanaey Igorots from the highlands of the Philippine Mountain Province as likely the  
9 closest living representatives of the source population that may have given rise to the  
10 Austronesian expansion. This conclusion rests on independent evidence from various analyses  
11 of autosomal data and uniparental markers.

12 Given the extensive presence of trade goods, cultural and linguistic evidence of Indian  
13 influence in Southeast Asia starting from 2.5kya we also detect traces of a South Asian  
14 signature in different populations in the region dating to the last couple of thousand years.

15

16 **Keywords:** population genetics, Southeast Asia, India, Austronesian expansion, Kankanaey

17

18

19

20

21

22

23

24

25

26

# 1 INTRODUCTION

2 Mainland (MSEA) and Island Southeast Asia (ISEA) are home to hundreds of different ethno-  
3 linguistic groups each displaying a complex demographic history.<sup>1</sup> Previous studies have  
4 revealed strong genetic correlations between populations which are geographically and  
5 linguistically close and suggested a common origin of all Southeast Asian and East Asian  
6 populations from a single migration wave.<sup>2</sup> It is well known, however, that in the more recent  
7 past the populations living in this region have undergone major demographic changes,  
8 particularly during the last five thousand years in association with the spread of the Neolithic  
9 cultural complex and Austronesian languages.<sup>3</sup> Wollstein and colleagues<sup>4</sup> reported significant  
10 genetic contributions from people currently inhabiting the Borneo (used as a proxy for Asian  
11 influence) and Papua New Guinea islands into Malayo-Polynesians (Austronesians who  
12 migrated beyond Taiwan) from Near and Remote Oceania. These admixture events were  
13 dated to approximately 3 kya, consistently with similar population movements involving  
14 people of Asian ancestry moving through ISEA dated around 4-3 kya.<sup>5</sup> More recent studies<sup>6,7</sup>  
15 have distinguished at least three major ancestral components in MSEA and ISEA in  
16 association with Papuan, Austro-Asiatic and Austronesian speaking populations. However the  
17 analyses aiming to identify the likely source regions of these dispersals are confounded by  
18 recent admixture in most modern ISEA populations with groups originating from other  
19 regions including MSEA<sup>2,8</sup> (see Text S1 for more details on the candidate populations  
20 included in this study).

21 In addition to the migratory events involving South East Asian sources, more recent South  
22 Asian influences in forms of cultural and trading networks, starting more than 2kya, in ISEA  
23 and MSEA have been well established from historical and archaeological data.<sup>9-12</sup>

24 Exemplary for these developments are sites the sites of Khao Sam Kaeo and Phu Khao Thong  
25 from Peninsular Thailand yielding archaeological evidence dating to 2.3-1.2kya.. They  
26 confirm the earliest trade networks with India, which include rouletted ware, semi-precious

1 stone beads and artefacts, and Indian crops.<sup>13</sup> In ISEA, one finds evidence of Indian trade  
2 either directly or via peninsular Thailand. Coastal sites located in Northern Bali dating to 2.1  
3 kya yielded pottery of East Indian or Sri Lankan production, gold and carnelian objects from  
4 North India and mung bean.<sup>14</sup> Furthermore epigraphy indicates a strong Indian impact on the  
5 nascent political structures of the region<sup>15</sup> and provides records of Brahmanic rituals and  
6 animal sacrifices<sup>16</sup>.

7 Linguistic evidence also supports early interethnic contact between Indian and Southeast  
8 Asian populations. Apart from the ubiquitous influence of Sanskrit<sup>17</sup> where it is difficult to  
9 distinguish ancient from more recent borrowings, analyses of the earliest Maritime Southeast  
10 Asian literature demonstrate that it already exhibits signs of Tamil influence from South  
11 India, much of which most likely spread across the region through pre-existing local  
12 networks.<sup>18</sup> Traces of paternal (Y chromosomes) and maternal (mtDNA) Indian ancestry have  
13 been detected across several Indonesian islands at low frequency (<5%).<sup>19-22</sup> The influx of  
14 Indian ancestry is detectable in some genome-wide analyses of low density autosomal SNP  
15 data<sup>2</sup> while being restricted to just a few populations from western Indonesia (Sumatra).  
16 Contrarily to that, a more recent study<sup>23</sup> using medium density SNP data could not find a  
17 South Asian genetic signature in South East Asia. The same authors however inferred gene  
18 flow from the Indian sub-continent to Aboriginal Australian populations and dated it at  
19 around 4kya. In the absence of a similar South Asian component in SEA this finding was  
20 interpreted to require a direct sea route bypassing Southeast Asia to explain such a signature  
21 in Australasia.

22 In order to refine the current understanding on the source of the Austronesian expansion and  
23 to further explore potential South Asian genetic contributions in MSEA and ISEA, we  
24 generated high density (730K) SNP Chip data for a panel of 196 individuals from 10  
25 populations including 50 of which (from the Bajo and Lebbo populations) are published  
26 already<sup>7</sup> and 146 new (Burmese and Vietnamese from MSEA, Ilocano, Tagalog and

1 Kankanaey from the Philippines, Murut, Malay and Dusun from ISEA plus 4 Australian  
2 Aborigines). We merged the newly generated dataset with those available from the literature  
3 (cf. Material and Methods) and i) investigated the existence of signs of South Asian admixture  
4 in our new SEA populations, ii) refined current knowledge on the putative source of the  
5 Austronesian expansion; iii) determined the extent to which signs of local adaptation are  
6 shared across local populations, as function of their common demographic history.

7

## 8 **MATERIAL AND METHODS**

### 9 **Samples, genotyping and phasing**

10 The newly generated dataset for this study consists of 150 individuals from 9 Southeast Asian  
11 and one Australian population (Figure 1 and Table S10). DNA was extracted from saliva  
12 samples collected from healthy adult donors who signed an informed consent form. The study  
13 was approved by local Research Ethics Committees (SingHealth Centralised Institutional  
14 Review Board and the Medical and Health Research Ethics Committee of the National Cancer  
15 Centre, Brunei Darussalam), the Cambridge Ethics Committee (HBREC.2011.01) and the  
16 ERC Ethics Panel. Southeast Asian samples were genotyped using Illumina OmniExpress  
17 Bead Chips for 730 525 SNPs. They are accessible together with 50 Bajo and Lebbo samples  
18 under the GEO accession number GSE77508. For the three Australian samples the Illumina  
19 Human 660K Quad Bead Chip yielded 655 215 SNPs, while for one Australian the 610K  
20 version of the latter chip gave 616 795 variants. These four samples are accessible under the  
21 accession number EGAS00001001738 in the European Genome-phenome Archive.

22 Before the analyses as such data filtering and quality checks using PLINK 1.07<sup>24</sup> were  
23 performed. Firstly, only autosomal SNPs with a genotyping success rate greater than 98%  
24 were included. PLINK was also utilized to remove all individuals more closely related than  
25 first degree cousins. This was done by estimating pairwise identity by descent (IBD)

1 iteratively; individuals with an IBD  $> 0.125$  were excluded. Following these quality controls  
2 haplotypes were inferred from genotype data with SHAPEIT.<sup>25</sup>  
3 Furthermore, 8 full mitochondrial Kankanaey genomes were sequenced by Complete  
4 Genomics (Mountain View, California, USA) using CG software version 2.4. Access to the  
5 sequences is provided under the GenBank accession numbers KU752558 to KU752565. All  
6 novel data from this paper will also be available under [www.ebc.ee/free\\_data](http://www.ebc.ee/free_data) in the PLINK  
7 (genotype data) and fasta (mitochondrial genomes) formats respectively.

8

### 9 **Demographic analyses**

10 To get a first overview for the novel Southeast Asian data we merged them with four  
11 reference populations from the HapMap 3 panel<sup>26</sup> and the HGDP Papuans<sup>27</sup> to obtain a set of  
12 307 625 common SNPs. Runs of Homozygosity (rOH), average observed heterozygosity and  
13 IBD were obtained using PLINK default parameters. Furthermore pairwise  $F_{ST}$  was calculated  
14 using an *ad hoc* Perl script implementing an estimator for Wright's formula.<sup>28</sup>

15 To address more specific questions regarding the ancestries of our novel populations we  
16 performed two distinct ADMIXTURE<sup>29,30</sup> analyses. For comparative purposes publicly  
17 available genotype data from the HapMap<sup>26</sup>, HDGP<sup>27</sup> and the Pan-Asian Consortium<sup>2</sup> projects  
18 was added to 185 individuals from 9 SEA populations (the divergence from the original  
19 number of 196 is due to the removal of close relatives). Additionally we used SNP data from  
20 studies focused on Indian populations.<sup>31,32</sup> This resulted in a dataset consisting of 1099  
21 individuals.

22 For further verification of our ADMIXTURE analysis, we assembled a second panel of 1010  
23 samples including 187 samples from our 9 SEA populations, and 4 Australian Aborigines,  
24 which are newly reported here. The samples, populations and references for both analyses are  
25 listed in Table S10. A detailed description of the merging and data curation for ADMIXTURE  
26 can be found in the Text S2.

1 Effective population size for our 9 SEA populations was estimated by analyzing LD patterns  
2 with the NeON R package.<sup>33</sup> To further investigate genetic structure and gene flow between  
3 populations we used the TreeMix v1.1 software package.<sup>34</sup> To measure how well the trees  
4 with different numbers of migration events (N) reflect the relationship between population  
5 groups we calculated the fraction  $f$  of explained variance as described by the original authors  
6 of the method. We used MEGA v6.0.6<sup>35</sup> to create a graphic representation of the TreeMix  
7 output. For specific admixture events of interest suggested by the ADMIXTURE plots the  
8 respective sets of recipient and source populations were tested with the three populations test  
9 ( $f_3$ ).<sup>34,36</sup> The population trios yielding a Z-score smaller than -2 were considered significantly  
10 admixed. These were then analyzed with ALDER<sup>37</sup> to date the putative admixture event.  
11 Furthermore we used the  $f_4$ -ratio test<sup>38</sup> to obtain a quantitative estimate of admixture  
12 percentages of interest.

13 For the analysis of the mtDNA data the haplogroup affiliation of each sample was assigned  
14 using HaploGrep 2.0<sup>39</sup> and PhyloTree build 16 (as of 19/02/2014)  
15 (<http://www.phylotree.org>).<sup>40</sup> The variants are described relative to the rCRS (GenBank  
16 Accession Number NC\_012920.1).<sup>41</sup>

17

18

## 19 **Selection tests**

20 To capture haplotype homozygosity based signals the Integrated Haplotype Score (iHS)<sup>42</sup> and  
21 Cross Population Extended Haplotype Homozygosity (XP-EHH)<sup>43</sup> tests were used. Both the  
22 iHS and XP-EHH statistics were calculated as in Pickrell et al. (2009)<sup>44</sup>, yielding about 10  
23 000-11 000 genomic windows for iHS and about 13,700 windows for XP-EHH for each SEA  
24 population analyzed. From the top 1% of all iHS signals, putatively the strongest candidates  
25 for selection, windows present in the top 5% iHS windows of the CHB population from the  
26 Hap Map panel were excluded, to pick up only signals particular to SEA. However, for the

1 analysis of regional sharing patterns based on the iHS this condition did not apply. For the  
2 XP-EHH the use of a reference population is inherent in the method, again CHB was chosen,  
3 for similar reasons.  
4 Furthermore, we computed the allele frequency based Population Branch Statistic (PBS). This  
5 test statistic represents the amount of allele frequency change at a given locus in the history of  
6 the test population since it diverged from other populations.<sup>45</sup> The outgroups for each tested  
7 SEA population were the YRI and CHB populations. Pairwise  $F_{ST}$  values for the populations  
8 of interest and the references were calculated following Weir and Cockerham.<sup>46</sup> PBS scores  
9 were estimated from the pairwise  $F_{ST}$  values.<sup>45</sup> Based on the approach of Pickrell et al.  
10 (2009)<sup>44</sup> the genome was divided into windows of a modified size of 100kb and the  
11 maximum PBS score in each window was used as the test statistic. This resulted in between  
12 26 000 and 27 000 windows for each analyzed group.

13

## 14 **RESULTS**

15 To investigate general patterns of population structure in our data we performed two distinct  
16 ADMIXTURE analyses: the first was mainly focused on populations from Southeast Asia and  
17 South Asia while the second provided the context of a broader, worldwide genetic landscape  
18 and additional validation for inferences from the first analysis.

19 According to the cross-validation scores for both analyses  $K=9$  admixture fractions provide  
20 the best fit (for the local plot additional  $K$ s are provided in Figure S2, for the global plot  $K$ s  
21 from 3 to 15 are shown in Figure S3B). The ADMIXTURE analyses of the newly generated  
22 data (Figure 1A, Figure S1) recapitulate the main ancestral components associated with  
23 Austronesian ( $k_6$ ), Austro-Asiatic ( $k_5$ ) and Papuan ( $k_3$ ) populations (Figure 1, Figure S2)  
24 already described in the area by previous studies.<sup>5,6</sup> At lower  $K$  values the component  
25 associated with the Papuans is highly prevalent in Eastern Indonesia and the Mamanwa (a



1 Negrito group from the Philippines) while at higher values it continues to persist only in the  
2 Alorese and Bajo from Indonesia (Figure 1B, Figure S2).

3 Burmese and Vietnamese exhibit significant proportions of the k2 component indicating  
4 shared ancestry with East Asian populations. The k4 component associated with South Asian  
5 ancestry is also consistently visible in Burmese and Malays (this study) and some Indonesian  
6 populations, mainly the Batak of Sumatra.<sup>2</sup> However at lower Ks this component is also  
7 present in the Javanese and the Mamanwa Negritos, suggesting affinities which however  
8 decline with higher Ks (Figure 1B, Figure S2). Notably, in the extended worldwide analysis  
9 (Figure S3B) the Papuan-related component (red) in the Bajo and the South Asian signal  
10 (green) in the Burmese and Malays were also clearly detectable. The SEA groups described  
11 here exhibit a remarkable diversity from very heterogeneous groups such as the Malays to the  
12 Kankanaey who appear homogenous in their ancestry composition by the ADMIXTURE  
13 analyses (Figure 1B, Figure S2).

14 The Kankanaey are an indigenous population of northern Luzon, belonging to the broader  
15 “Igorot” group. At K=9, the majority of Kankanaey ancestry is in the k6 component, which  
16 they share with the Ami (AX-AM) and Atayal (AX-AT) from Taiwan and, hence, is  
17 putatively associated with the Austronesian expansion (Figure 1A, Figure S2). When it  
18 emerges as distinct from the other Asian components, the k6 brown ancestry is spread  
19 throughout ISEA and remains stable in all these groups from K8-10 (Figure 1B, Figure S2).

20 Remarkably in the regional admixture plots the Kankanaey remain unadmixed throughout all  
21 Ks from 2-10 (Figure S2), even though at lower Ks they do not yet have their own distinct  
22 component. These findings are consistent with the Kankanaey’s geographic location, the  
23 Mountain Province in the Northern Philippines (Figure 1A, Figure S1), close to Taiwan, the  
24 likely center of the Austronesian expansion.<sup>3,6</sup>

25 Kankanaey genome wide heterozygosity levels and extent of runs of homozygosity (rOH)  
26 (Table S1) rule out potential confounders such as extreme inbreeding or genetic drift being

1 causative for their unusually homogeneous ancestry. To further explore the potential effect of  
2 demographic history on population structure we estimated the effective population size of the  
3 nine SEA populations presented here based on the development of linkage disequilibrium  
4 (LD) patterns over time (Figure S4).<sup>33</sup> The mainland Burmese and Vietnamese groups exhibit  
5 comparatively high effective population sizes and signs of recent expansion. This is in line  
6 with their recent history of admixture with neighboring populations, whereas there is more  
7 variation in the ISEA populations. Notably the Kankanaey have one of the lowest values  
8 varying between 2000-3000 (6000-27000 kya). However they are not an extreme outlier and  
9 are comparable to the Lebbo from Borneo (no significant difference,  $p = 0.7938$ ), who instead  
10 do not show such a homogeneous ADMIXTURE profile. Under the assumption that the  
11 brown k6 component reflects ancestry connected to the Austronesian expansion, the  
12 Kankanaey displayed a higher percentage of it than even Austronesian Taiwanese populations  
13 (AX-AT, AX-AM, Figure 1A, Figure S2). The affinity of the Kankanaey to these groups was  
14 supported by the TreeMix<sup>34</sup> analyses of 25 populations (Figures S5-S6) where the Kankanaey  
15 did not cluster with other Filipinos but rather with the Taiwanese aboriginals.

16 The emerging picture seems to be compatible with a scenario of local Austroasiatic and  
17 Papuan components influenced by the incoming Austronesian (brown k6, Figure 1A, Figure  
18 S2) wave 4-3kya which originated from a population living in Taiwan and, perhaps, in the  
19 North Philippines.<sup>6</sup> The attempt to date the above admixture events using ALDER<sup>37</sup>  
20 highlighted a clear admixture pattern between “Kankanaey like” people and earlier substrates,  
21 dated to at least 2.2 kya in the Bajo (Table 1).

22 These affinities of the Kankanaey and their potential role as a good proxy for the  
23 Austronesian expansion are further highlighted when looking at uniparental markers. The  
24 eight available Kankanaey mtDNA sequences (Table S2) exhibit lineages (B4a1a;M7b1a2a1)  
25 which are typical markers of Malayo-Polynesian speaking populations.<sup>47,48</sup>

1 Lastly, the Kankanaey cannot be modeled as any kind of mixture from 46 populations using  
2 the  $f_3$  statistic (Table S3).<sup>36</sup> Taken together, the evidence from these independent approaches  
3 suggests that the Kankanaey could potentially represent an unadmixed remnant population  
4 close to the source that may have given rise to the Austronesian expansion.

5 We also utilized the  $f_3$  test together with ALDER to further contextualize the potential South  
6 Asian connections of some SEA groups. Both of these statistics (Table 1) suggest the  
7 presence of variable degrees of South Asian-related ancestry in the MSEA and ISEA  
8 populations (Bajo, Burmese, Filipino and Malay). Assuming a generation time of 30 years<sup>49</sup>  
9 the earliest possible midpoint of the South Asian admixture is estimated at 2.4 kya. The  
10 overall proportion of South Asian ancestry was further estimated by applying the  $f_4$  statistic<sup>38</sup>  
11 (Table S4) according to the tree presented in Figure S7. The estimated values were 24.9% for  
12 the Burmese, 8.3% for the Malays and 5.3% for the Bajo. One limitation of this approach is  
13 its dependence on shared genetic drift. As the Papuans and South Indians have a similar  
14 position in the phylogenetic tree relative to the other groups, Papuan ancestry could be  
15 mistaken as South Indian. This has probably no effect in the Burmese and Malay, who do not  
16 show Papuan admixture (Figure 1A, Figure S2) but could contribute to the South Indian  
17 ancestry detected in the Bajo. True Indian ancestry in this population still seems conceivable  
18 given the presence of South Asian lineages in uniparental marker analyses<sup>22</sup>

19 These analyses indicate a South-Asian related component in the genetic make-up of at least  
20 some SEA groups which entered their gene pool ca. 2.4 kya ago, being supported by  
21 ADMIXTURE,  $f_3$  and  $f_4$  analyses for the Burmese and the Malay and by  $f$ -statistics for the  
22 Bajo ( $f_3$ ,  $f_4$ ) and the lowland Filipinos ( $f_3$ ).

23 As an additional tool to explore relationships among populations, we examined patterns of  
24 haplotype homozygosity and allelic differentiation using test statistics  $iHS$ <sup>42</sup>,  $XP-EHH$ <sup>43</sup>, and  
25  $PBS$  test<sup>45</sup> (Tables S7, S8, S9). For the  $iHS$  the amount of signal sharing between two groups  
26 correlates only very weakly ( $r^2=0.041$  for a linear regression) to overall genetic similarity as

1 expressed by the  $F_{ST}$  (Figure S8). However, the MSEA groups and the Han Chinese (included  
2 as a reference) who share a considerable proportion of East-Asian ancestry (Figure 1A, Figure  
3 S2) also show a great affinity to each other regarding haplotype homozygosity patterns (Table  
4 S5). In ISEA those groups with at least three significant ancestry components at  $K=9$  (Bajo,  
5 Filipino, Malay, Figure 1A) exhibit more signal sharing. In contrast, Kakanaey, Lebbo and  
6 Murut show reduced sharing with all other populations, which perhaps highlights phenomena  
7 of deep population splits and separate demographic histories in recent times when the  
8 haplotype homozygosities have accumulated.

9 However, these inferences are highly dependent upon the approach utilized. A different  
10 picture presents itself for the XP-EHH, which considers both haplotype homozygosity and  
11 allelic differentiation, with the Han Chinese used as outgroup. The average fraction of signal  
12 sharing declines from 0.31 to 0.22, while the correlation with the  $F_{ST}$  increases considerably  
13 ( $r^2= 0.256$ ). This is probably because signals connected to shared ancestry with East Asians  
14 are excluded. It causes the Burmese, who exhibit a large fraction of the  $k2$  East Asian-related  
15 component (Figure 1A, Figure S2) to become an outlier especially with respect to their high  
16 fraction of unique top 1% XP-EHH signals, only 15% of which are shared with other groups  
17 on average.

18

## 19 **DISCUSSION**

20 In this study we set out to explore population structure in MSEA and ISEA and more  
21 specifically, to clarify the exact nature of South Asian gene flow into SEA and the presence of  
22 potential un-admixed Austronesian population(s) close to the ancestral Austronesian source.

23 We detected a minor South Asian component in our ADMIXTURE analyses in MSEA and  
24 ISEA populations (green  $k4$ , Figure 1A, Figure S2; green, Figure S3B) which was further  
25 confirmed by  $f3$ ,  $f4$  and ALDER results and dated to have entered SEA from 2.4kya (Table  
26 1). While this component is more widespread at lower  $Ks$  (Figure 1B, Figure S2) at the best

1 K=9 (Figure 1A) the evidence is strongest for the Burmese and the Malay and somewhat  
2 weaker for Bajo and Filipinos, where it is limited to the f statistics (Table 1, Table S4). It is  
3 important to explore how these results relate to the linguistic and archaeological evidences,  
4 attesting a continuous presence of South Asian cultures in Southeast Asia since 2.5  
5 kya.<sup>12,17,50,51</sup> This should be done keeping in mind that in the majority of SEA populations the  
6 Indian component is absent or below the scale of a potential error and detectability. Firstly, it  
7 is most likely that the “carriers” of South Asian culture were traders, artisans<sup>50</sup> and at a later  
8 date, religious scholars (Brahmins) who were influential as advisers to Southeast Asian rulers.  
9 Some of these might have been locals educated in India who brought home Sanskrit texts and  
10 Brahmanic rituals.<sup>52</sup> So this rather small group would not have left a major genetic signature.  
11 Secondly, the epigraphic record and evidence from monumental archaeology during the late  
12 first millennium CE attests that the Indian presence is biased towards courts and generally  
13 higher social strata, which can lead us to overestimate the impact on the majority of the  
14 population.<sup>52</sup> More generally speaking there are a wide range of scenarios relating to the  
15 spread of cultural elements and gene flow and the patterns of this relationship are highly  
16 complex to model (cf. the example of the Neolithization in Europe<sup>53</sup>). So with the exception  
17 of the Burmese, who are also geographically very close to the Indian subcontinent, the  
18 evidence points to rather minor Indian gene flow, in contrast to the documented cultural  
19 influence which, however, overlaps with the admixture range dated with ALDER (Table 1).  
20 This low South Asian gene flow was however also detected in some other populations across  
21 ISEA.<sup>2,19-22</sup> Taken together these findings suggest Southeast Asia as a potential waypoint for  
22 the reported South Asian migration into Australasia which was disputed by the authors who  
23 proposed this migration event.<sup>23</sup> However the date obtained using ALDER (2.4 kya) is at least  
24 1500 years posterior to the reported South Indian migration into Australasia.<sup>23</sup> A preliminary  
25 conclusion would envisage the SEA and Australasia migrations as two separate events.  
26 Besides the fact that the dating methods were different in our case and Pugach et al. (they

1 used a method based on wavelet transform analysis) at least two caveats can be brought up to  
2 reconcile this fragmented scenario. Given the evidence presented here, it seems reasonable to  
3 assume a constant gene flow from South Asia into SEA via land, with Australasia being only  
4 a sporadic end-point. In this case the 4kya estimate provided by Pugach and colleagues would  
5 be a point estimate of the sparse arrival into Australasia, while our ALDER estimate should  
6 be interpreted as the midpoint<sup>37</sup> of such a flow between 4kya and more recent times.  
7 Secondly, given the surprising concordance of linguistic and archaeological evidences for a  
8 South Asian presence in SEA around 2.5 kya, one could imagine a particularly intense  
9 corresponding gene flow during that time further biasing the ALDER estimate toward this  
10 period.

11 In this study we have identified the Kankanaey from the northern Philippines as the  
12 population harboring the highest reported amount of the Austronesian genomic component,  
13 even higher than the ones detectable in modern aboriginal Taiwanese (Figure 1B, Figure S2).  
14 This conclusion rests on evidence from several independent analyses including  
15 ADMIXTURE, f3, rOH, TreeMix,  $N_e$  and uniparental markers.

16 The Kankanaey belong to the broader group of populations collectively known as Igorot  
17 (Text S1). Various studies exist on the Kankanaey language<sup>54</sup> and customs<sup>55</sup>, although works  
18 on their prehistory are lacking. Genetic data from 30 Kankanaey-speakers was included in a  
19 recent study of the mtDNA-haplotype-diversity in the Philippines.<sup>56</sup> There they were shown to  
20 share many lineages with two other Igorot groups (Ibaloi and Ifugao) from Northern Luzon.  
21 These results are broadly consistent with the uniparental data we present here (Table S2),  
22 where the Kankanaey show haplotypes also found in Taiwanese aboriginals<sup>57</sup> and generally  
23 associated with the Austronesian expansion<sup>47,48</sup>. We conclude that the Kankanaey are either  
24 the best preserved source of the Austronesian expansion, or a case of total replacement that  
25 followed it. The dominant model suggests a southward diffusion of Austronesians from  
26 Taiwan around 4000 BP, which impacted the Philippines, the north of Borneo and Sulawesi

1 between 3800–3600 BP, and later spread into the Pacific.<sup>3</sup> Even if the modality of this  
2 expansion is complex and still debated<sup>58</sup>, the location of the Kankanaey in the northern  
3 Philippines, close to Taiwan, suggests that they may be considered as one of the least admixed  
4 living groups tracing their ancestry from the source populations of the Austronesian  
5 expansion. Furthermore we confirm the finding of an Austro-Asiatic-related component in  
6 ISEA populations (here the Dusun, Murut, Lebbo and Bajo) first reported by Lipson et al. in  
7 2014<sup>6</sup> and there described as unexpected due to the historically nearly exclusive presence of  
8 Austro-Asiatic speakers on the mainland. Given its wide spread in MSEA and ISEA in  
9 linguistically diverse groups, the explicit association of k5 with this language family should  
10 be taken with caution. However, it is worthwhile noting that in India we find this component  
11 specifically in Munda speaking populations. The k5 component could represent an ancestral  
12 substrate, which was once distributed widely throughout SEA and was encountered by the  
13 Austronesians when they spread from Taiwan. Another possibility is that there was an early  
14 split into several subgroups during the Austronesian expansion and that this component  
15 belongs to the ancestral make-up of a subgroup of Malayo-Polynesians who expanded into  
16 western Indonesia.

17 Our comparison of haplotype-based scans of positive selection revealed that compared to  
18 earlier studies on a continental level<sup>32</sup> in a regional context in ISEA there is no good  
19 correlation between haplotype sharing patterns and genetic distance as indicated by the  $F_{ST}$   
20 (Figure S8). However, as described above, the haplotype homozygosity patterns still reflect  
21 demography to a considerable extent. Populations showing more diversity in the admixture  
22 plots also exhibit higher levels of shared signals with other groups. Furthermore, the sharing  
23 patterns proved to be very dependent on the kind of test utilized. Notably when the XP-EHH,  
24 which uses the Han Chinese as outgroup, is applied, all signals shared with East Asians are  
25 excluded. Intriguingly this causes the Burmese whose ancestry contains a significant South

1 Asian-related component (Figure 1A, Figure S2) to become an outlier (Table S6) potentially  
2 reflecting haplotype homozygosity signals unique to their share of Indian ancestry.  
3 In conclusion, we report a minor South Asian contribution to the genomes of some modern  
4 MSEA and ISEA populations, mainly the Burmese and the Malay. This is in line with a  
5 general cultural diffusion process to SEA, driven by smaller groups of influential individuals  
6 from South Asia. Secondly, our work strongly suggests that based on the currently available  
7 data the Kankanaey tribal group from Northern Luzon, Philippines are the best genetic  
8 representative of the Austronesian expansion. We envisage high coverage whole genome  
9 sequencing of this population as a sound approach to further explore this major peopling  
10 event that shaped the genetic landscape of the broader South East Asia region.

11  
12

### 13 **Acknowledgements**

14 This work was supported by the European Research Council Starting Investigator grant FP7-  
15 261213 to T.K. This work was supported by the French ANR grant number ANR-14-CE31-  
16 0013-01 (grant OceoAdapto to F-X.R.), the French ANR-12-PDOC-0037-01 (grant  
17 GENOMIX to D.P.), the Region Aquitaine of France (grant MAGE to T.L.), the French  
18 Ministry of Foreign and European Affairs (French Archaeological Mission in Borneo  
19 (MAFBO) to F-X.R). T.A. was supported by a Wellcome Trust Post-Doctoral fellowship  
20 (WT1000MA).

21

### 22 **CONFLICT OF INTEREST**

23

24 The authors declare no conflict of interest.

25

26 Supplementary Information is available at the European Journal of Human Genetics' website.

27



## 1 REFERENCES

- 2 1 Lewis MP, Simons GF, Fennig CD (eds.). *Ethnologue: Languages of the World*,  
3 *Eighteenth edition*. SIL International: Dallas, Texas, 2015<http://www.ethnologue.com>.
- 4 2 HUGO Pan-Asian SNP Consortium, Abdulla MA, Ahmed I, Assawamakin A, Bhak J,  
5 Brahmachari SK *et al*. Mapping human genetic diversity in Asia. *Science* 2009; **326**:  
6 1541–1545.
- 7 3 Bellwood PS. *Prehistory of the Indo-Malaysian Archipelago*. ANU E Press: Canberra,  
8 2007<http://epress.anu.edu.au/?p=80041> (accessed 26 Nov2015).
- 9 4 Wollstein A, Lao O, Becker C, Brauer S, Trent RJ, Nürnberg P *et al*. Demographic  
10 history of Oceania inferred from genome-wide data. *Curr Biol CB* 2010; **20**: 1983–1992.
- 11 5 Xu S, Pugach I, Stoneking M, Kayser M, Jin L. Genetic dating indicates that the Asian–  
12 Papuan admixture through Eastern Indonesia corresponds to the Austronesian expansion.  
13 *Proc Natl Acad Sci* 2012; : 201118892.
- 14 6 Lipson M, Loh P-R, Patterson N, Moorjani P, Ko Y-C, Stoneking M *et al*. Reconstructing  
15 Austronesian population history in Island Southeast Asia. *Nat Commun* 2014; **5**: 4689.
- 16 7 Pierron D, Razafindrazaka H, Pagani L, Ricaut F-X, Antao T, Capredon M *et al*.  
17 Genome-wide evidence of Austronesian-Bantu admixture and cultural reversion in a  
18 hunter-gatherer group of Madagascar. *Proc Natl Acad Sci* 2014; **111**: 936–941.
- 19 8 Trejaut JA, Poloni ES, Yen J-C, Lai Y-H, Loo J-H, Lee C-L *et al*. Taiwan Y-  
20 chromosomal DNA variation and its relationship with Island Southeast Asia. *BMC Genet*  
21 2014; **15**: 77.
- 22 9 Ardika W, Bellwood PS. Sembiran: the beginnings of Indian contact with Bali. *Antiquity*  
23 1991; **65**: 221–232.
- 24 10 Ardika W, Bellwood PS, Sutaba IM, Yuliati KC. Sembiran and the first Indian contacts  
25 with Bali: an update. *Antiquity* 1997; **71**: 193–95.
- 26 11 Lawler A. Sailing Sinbad’s seas. *Science* 2014; **344**: 1440–1445.

- 1 12 Manguin P-Y, Mani A, Wade G (eds.). *Early interactions between South and Southeast*  
2 *Asia: reflections on cross-cultural exchange*. Institute of Southeast Asian Studies ;  
3 Manohar India: Singapore : New Delhi, 2011.
- 4 13 Castillo C. *The Archaeobotany of Khao Sam Kaeo and Phu Khao Thong: The Agriculture*  
5 *of Late Prehistoric Southern Thailand*. 2013.
- 6 14 Calo A, Prasetyo B, Bellwood P, Lankton JW, Gratuze B, Pryce TO *et al*. Sembiran and  
7 Pacung on the north coast of Bali: a strategic crossroads for early trans-Asiatic exchange.  
8 *Antiquity* 2015; **89**: 378–396.
- 9 15 Mabbett IW. The ‘Indianization’ of Southeast Asia: Reflections on the Historical Sources.  
10 *J Southeast Asian Stud* 1977; **8**: 143–161.
- 11 16 Guy J. Tamil merchants and the Hindu-Buddhist Diaspora in early Southeast Asia. In:  
12 Manguin P-Y, Mani A, Wade G (eds). *Early interactions between South and Southeast*  
13 *Asia: reflections on cross-cultural exchange*. Institute of Southeast Asian Studies ;  
14 Manohar India: Singapore : New Delhi, 2011, pp 243–262.
- 15 17 Gonda J. *Sanskrit in Indonesia*. International Academy of Indian Culture.: New Dehli,  
16 1973.
- 17 18 Hoogervorst T. Detecting pre-modern lexical influence from South India in Maritime  
18 Southeast Asia. *Archipel* 2015; **89**: 63–93.
- 19 19 Chaubey G, Endicott P. The Andaman Islanders in a regional genetic context:  
20 reexamining the evidence for an early peopling of the archipelago from South Asia. *Hum*  
21 *Biol* 2013; **85**: 153–172.
- 22 20 Karafet TM, Lansing JS, Redd AJ, Reznikova S, Watkins JC, Surata SPK *et al*. Balinese  
23 Y-chromosome perspective on the peopling of Indonesia: genetic contributions from pre-  
24 neolithic hunter-gatherers, Austronesian farmers, and Indian traders. *Hum Biol* 2005; **77**:  
25 93–114.

- 1 21 Karafet TM, Hallmark B, Cox MP, Sudoyo H, Downey S, Lansing JS *et al.* Major East-  
2 West Division Underlies Y Chromosome Stratification across Indonesia. *Mol Biol Evol*  
3 2010; **27**: 1833–1844.
- 4 22 Kusuma P, Cox MP, Brucato N, Sudoyo H, Letellier T, Ricaut F-X. Western Eurasian  
5 genetic influences in the Indonesian archipelago. *Quat Int* 2015.  
6 doi:10.1016/j.quaint.2015.06.048.
- 7 23 Pugach I, Delfin F, Gunnarsdóttir E, Kayser M, Stoneking M. Genome-wide data  
8 substantiate Holocene gene flow from India to Australia. *Proc Natl Acad Sci* 2013; :  
9 201211927.
- 10 24 Purcell S, Neale B, Todd-Brown K, Thomas L, Ferreira MAR, Bender D *et al.* PLINK: A  
11 Tool Set for Whole-Genome Association and Population-Based Linkage Analyses. *Am J*  
12 *Hum Genet* 2007; **81**: 559–575.
- 13 25 Delaneau O, Zagury J-F, Marchini J. Improved whole-chromosome phasing for disease  
14 and population genetic studies. *Nat Methods* 2013; **10**: 5–6.
- 15 26 International HapMap Consortium, Frazer KA, Ballinger DG, Cox DR, Hinds DA, Stuve  
16 LL *et al.* A second generation human haplotype map of over 3.1 million SNPs. *Nature*  
17 2007; **449**: 851–861.
- 18 27 Li JZ, Absher DM, Tang H, Southwick AM, Casto AM, Ramachandran S *et al.*  
19 Worldwide Human Relationships Inferred from Genome-Wide Patterns of Variation.  
20 *Science* 2008; **319**: 1100–1104.
- 21 28 Wright S. Evolution in Mendelian Populations. *Genetics* 1931; **16**: 97–159.
- 22 29 Alexander DH, Novembre J, Lange K. Fast model-based estimation of ancestry in  
23 unrelated individuals. *Genome Res* 2009; **19**: 1655–1664.
- 24 30 Cardona A, Pagani L, Antao T, Lawson DJ, Eichstaedt CA, Yngvadottir B *et al.* Genome-  
25 Wide Analysis of Cold Adaptation in Indigenous Siberian Populations. *PLoS ONE* 2014;  
26 **9**: e98076.

- 1 31 Chaubey G, Metspalu M, Choi Y, Magi R, Romero IG, Soares P *et al.* Population Genetic  
2 Structure in Indian Austroasiatic Speakers: The Role of Landscape Barriers and Sex-  
3 Specific Admixture. *Mol Biol Evol* 2011; **28**: 1013–1024.
- 4 32 Metspalu M, Romero IG, Yunusbayev B, Chaubey G, Mallick CB, Hudjashov G *et al.*  
5 Shared and Unique Components of Human Population Structure and Genome-Wide  
6 Signals of Positive Selection in South Asia. *Am J Hum Genet* 2011; **89**: 731–744.
- 7 33 Mezzavilla M, Ghirotto S. Neon: An R Package to Estimate Human Effective Population  
8 Size and Divergence Time from Patterns of Linkage Disequilibrium between SNPS. *J*  
9 *Comput Sci Syst Biol* 2015; **8**. doi:10.4172/jcsb.1000168.
- 10 34 Pickrell JK, Pritchard JK. Inference of population splits and mixtures from genome-wide  
11 allele frequency data. *PLoS Genet* 2012; **8**: e1002967.
- 12 35 Tamura K, Stecher G, Peterson D, Filipinski A, Kumar S. MEGA6: Molecular Evolutionary  
13 Genetics Analysis Version 6.0. *Mol Biol Evol* 2013; **30**: 2725–2729.
- 14 36 Reich D, Thangaraj K, Patterson N, Price AL, Singh L. Reconstructing Indian population  
15 history. *Nature* 2009; **461**: 489–494.
- 16 37 Loh P-R, Lipson M, Patterson N, Moorjani P, Pickrell JK, Reich D *et al.* Inferring  
17 admixture histories of human populations using linkage disequilibrium. *Genetics* 2013;  
18 **193**: 1233–1254.
- 19 38 Moorjani P, Patterson N, Hirschhorn JN, Keinan A, Hao L, Atzmon G *et al.* The History  
20 of African Gene Flow into Southern Europeans, Levantines, and Jews. *PLoS Genet* 2011;  
21 **7**: e1001373.
- 22 39 Kloss-Brandstätter A, Pacher D, Schönherr S, Weissensteiner H, Binna R, Specht G *et al.*  
23 HaploGrep: a fast and reliable algorithm for automatic classification of mitochondrial  
24 DNA haplogroups. *Hum Mutat* 2011; **32**: 25–32.
- 25 40 van Oven M, Kayser M. Updated comprehensive phylogenetic tree of global human  
26 mitochondrial DNA variation. *Hum Mutat* 2009; **30**: E386–E394.

- 1 41 Andrews RM, Kubacka I, Chinnery PF, Lightowers RN, Turnbull DM, Howell N.  
2 Reanalysis and revision of the Cambridge reference sequence for human mitochondrial  
3 DNA. *Nat Genet* 1999; **23**: 147.
- 4 42 Voight BF, Kudravalli S, Wen X, Pritchard JK. A Map of Recent Positive Selection in  
5 the Human Genome. *PLoS Biol* 2006; **4**: e72.
- 6 43 Sabeti PC, Varilly P, Fry B, Lohmueller J, Hostetter E, Cotsapas C *et al.* Genome-wide  
7 detection and characterization of positive selection in human populations. *Nature* 2007;  
8 **449**: 913–918.
- 9 44 Pickrell JK, Coop G, Novembre J, Kudravalli S, Li JZ, Absher D *et al.* Signals of recent  
10 positive selection in a worldwide sample of human populations. *Genome Res* 2009.  
11 doi:10.1101/gr.087577.108.
- 12 45 Yi X, Liang Y, Huerta-Sanchez E, Jin X, Cuo ZXP, Pool JE *et al.* Sequencing of 50  
13 Human Exomes Reveals Adaptation to High Altitude. *Science* 2010; **329**: 75–78.
- 14 46 Weir BS, Cockerham CC. Estimating F-Statistics for the Analysis of Population  
15 Structure. *Evolution* 1984; **38**: 1358–1370.
- 16 47 Trejaut JA, Kivisild T, Loo JH, Lee CL, He CL, Hsu CJ *et al.* Traces of Archaic  
17 Mitochondrial Lineages Persist in Austronesian-Speaking Formosan Populations. *PLoS*  
18 *Biol* 2005; **3**: e247.
- 19 48 Soares P, Rito T, Trejaut J, Mormina M, Hill C, Tinkler-Hundal E *et al.* Ancient  
20 Voyaging and Polynesian Origins. *Am J Hum Genet* 2011; **88**: 239–247.
- 21 49 Fenner JN. Cross-cultural estimation of the human generation interval for use in genetics-  
22 based population divergence studies. *Am J Phys Anthropol* 2005; **128**: 415–423.
- 23 50 Bellina B, Silapanth P, Chaisuwan B, Thongcharoenchaikit C, Bernard V, Borell B *et al.*  
24 The Development of Coastal Polities in the Upper Thai-Malay Peninsula in the Late First  
25 Millennium BCE. In: Reire N, Murphy SA (eds). *Before Siam: Essays in Art and*  
26 *Archaeology*. River Books: Bangkok, 2014, pp 69–89.

- 1 51 Calo A. Ancient trade between India and Indonesia. *Science* 2014; **345**: 1255–1255.
- 2 52 Bronkhorst J. The Spread of Sanskrit in Southeast Asia. In: Manguin P-Y, Mani A, Wade  
3 G (eds). *Early interactions between South and Southeast Asia: reflections on cross-*  
4 *cultural exchange*. Institute of Southeast Asian Studies ; Manohar India: Singapore : New  
5 Delhi, 2011, pp 263–275.
- 6 53 Fort J. Demic and cultural diffusion propagated the Neolithic transition across different  
7 regions of Europe. *J R Soc Interface* 2015; **12**: 20150166–20150166.
- 8 54 Allen JL. *Kankanaey: a role and reference grammar analysis*. SIL International  
9 Publications: Dallas, Texas, 2014.
- 10 55 Kohnen N. ‘Natural’ childbirth among the Kankanaly-Igorot. *Bull N Y Acad Med* 1986;  
11 **62**: 768–777.
- 12 56 Delfin F, Min-Shan Ko A, Li M, Gunnarsdóttir ED, Tabbada KA, Salvador JM *et al*.  
13 Complete mtDNA genomes of Filipino ethnolinguistic groups: a melting pot of recent and  
14 ancient lineages in the Asia-Pacific region. *Eur J Hum Genet* 2014; **22**: 228–237.
- 15 57 Ko AM-S, Chen C-Y, Fu Q, Delfin F, Li M, Chiu H-L *et al*. Early Austronesians: Into  
16 and Out Of Taiwan. *Am J Hum Genet* 2014; **94**: 426–436.
- 17 58 Bulbeck F. An integrated perspective on the Austronesian Diaspora: the switch from  
18 cereal agriculture to maritime foraging in the colonisation of Island Southeast Asia. *Aust*  
19 *Archaeol* 2008; **67**: 31–52.

## 21 **Figure Legends**

22  
23 **Figure 1:**  
24 **(A) A map of Southeast Asia, displaying a subset of populations assessed in this study**  
25 **and the distribution of ancestry components based on the local ADMIXTURE run with**  
26 **the optimal number of ancestry components (K=9, cf. Figure S2). The figure legend on**  
27 **the lower left section shows the list of genetic ancestry components whose color codes**  
28 **correspond to those on the pie charts. Components k8 and k9 are mainly present in the**  
29 **Yoruba and Ati Negritos respectively and do not significantly contribute to the genetic**  
30 **diversity of the groups displayed in Figure 1.**

1 **The population abbreviations are as follows: Alo-Alorese, Baj-Bajo, Bat-Batak,**  
2 **Bru-Brunei (Dusun, Murut), Bur- Burmese, CHB-Chinese from Beijing, Jav-Javanese,**  
3 **Kan-Kankanaey Igorots, Leb-Lebbo, Mal-Malay, Mam-Mamanwa Negritos, Men-**  
4 **Mentawai, Mun-Mundari, NIn-North Indians, Pap-Papuans, PhU-Philippine Urban,**  
5 **SIn-South Indians, Taw-Ami and Atayal from Taiwan, Viet-Vietnamese.**  
6 **Note that the symbols next to the population names reflect the linguistic affiliations.**  
7 **Austro-Asiatic languages: circle, Austronesian languages: asterisk, Indo-European**  
8 **languages: square, Dravidian languages: hash, Papuan languages: cross, Tibeto-Burman**  
9 **languages: caret.**  
10 **(B) Three graphs showing the proportions of ancestry components k3, k4 and k6 from**  
11 **their emergence as independent components in the Papuans (k3, red), Indian**  
12 **populations (k4, green) and the Kankanaey Igorot (k6, brown) across multiple higher K**  
13 **values. All populations displayed show a percentage of at least 5% of the respective**  
14 **ancestry when it emerges.**