



**HAL**  
open science

## Contrasting Paternal and Maternal Genetic Histories of Thai and Lao Populations

Wibhu Kutanan, Jatupol Kampuansai, Metawee Srikumool, Andrea Brunelli, Silvia Ghirotto, Leonardo Arias, Enrico Macholdt, Alexander Hübner, Roland Schröder, Mark Stoneking

► **To cite this version:**

Wibhu Kutanan, Jatupol Kampuansai, Metawee Srikumool, Andrea Brunelli, Silvia Ghirotto, et al.. Contrasting Paternal and Maternal Genetic Histories of Thai and Lao Populations. *Molecular Biology and Evolution*, 2019, 36, pp.1490 - 1506. 10.1093/molbev/msz083 . hal-04604766

**HAL Id: hal-04604766**

**<https://cnrs.hal.science/hal-04604766>**

Submitted on 7 Jun 2024

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# Contrasting Paternal and Maternal Genetic Histories of Thai and Lao Populations

Wibhu Kutanan,<sup>\*,1,2</sup> Jatupol Kampaunsai,<sup>3,4</sup> Metawee Srikumool,<sup>5</sup> Andrea Brunelli,<sup>6</sup> Silvia Ghirotto,<sup>6</sup> Leonardo Arias,<sup>2</sup> Enrico Macholdt,<sup>2</sup> Alexander Hübner,<sup>2</sup> Roland Schröder,<sup>2</sup> and Mark Stoneking<sup>\*,2</sup>

<sup>1</sup>Department of Biology, Faculty of Science, Khon Kaen University, Khon Kaen, Thailand

<sup>2</sup>Department of Evolutionary Genetics, Max Planck Institute for Evolutionary Anthropology, Leipzig, Germany

<sup>3</sup>Department of Biology, Faculty of Science, Chiang Mai University, Chiang Mai, Thailand

<sup>4</sup>Center of Excellence in Bioresources for Agriculture, Industry and Medicine, Chiang Mai University, Chiang Mai, Thailand

<sup>5</sup>Department of Biochemistry, Faculty of Medical Science, Naresuan University, Phitsanulok, Thailand

<sup>6</sup>Department of Life Science and Biotechnology, University of Ferrara, Ferrara, Italy

\*Corresponding authors: E-mails: wibhu@kku.ac.th; stoneking@eva.mpg.de.

Associate editor: Connie Mulligan

## Abstract

The human demographic history of Mainland Southeast Asia (MSEA) has not been well studied; in particular, there have been very few sequence-based studies of variation in the male-specific portions of the Y chromosome (MSY). Here, we report new MSY sequences of ~2.3 mB from 914 males and combine these with previous data for a total of 928 MSY sequences belonging to 59 populations from Thailand and Laos who speak languages belonging to three major Mainland Southeast Asia families: Austroasiatic, Tai-Kadai, and Sino-Tibetan. Among the 92 MSY haplogroups, two main MSY lineages (O1b1a1a\* [O-M95\*] and O2a\* [O-M324\*]) contribute substantially to the paternal genetic makeup of Thailand and Laos. We also analyze complete mitochondrial DNA genome sequences published previously from the same groups and find contrasting pattern of male and female genetic variation and demographic expansions, especially for the hill tribes, Mon, and some major Thai groups. In particular, we detect an effect of postmarital residence pattern on genetic diversity in patrilocal versus matrilocal groups. Additionally, both male and female demographic expansions were observed during the early Mesolithic (~10 ka), with two later major male-specific expansions during the Neolithic period (~4–5 ka) and the Bronze/Iron Age (~2.0–2.5 ka). These two later expansions are characteristic of the modern Austroasiatic and Tai-Kadai groups, respectively, consistent with recent ancient DNA studies. We simulate MSY data based on three demographic models (continuous migration, demic diffusion, and cultural diffusion) of major Thai groups and find different results from mitochondrial DNA simulations, supporting contrasting male and female genetic histories.

**Key words:** Y chromosome, mtDNA, Austroasiatic, Tai-Kadai, Sino-Tibetan.

## Introduction

Thailand and Laos occupy a key location in the center of Mainland Southeast Asia (MSEA; [fig. 1](#)), which is undoubtedly one of the factors facilitating the extensive ethnolinguistic diversity, as there are 68 recognized groups in Thailand and 82 groups in Laos, belonging to five language families ([Simons and Fennig 2018](#)). The prehistoric peopling of the area of present-day Thailand and Laos has been documented by several archaeological studies ([Shoocongdej 2006](#); [Demeter et al. 2012](#); [Higham 2014, 2017](#)) and investigated further by recent ancient DNA studies ([Lipson et al. 2018](#); [McColl et al. 2018](#)). The earliest presence of modern humans in SEA is dated to ~50 ka ([Higham 2013](#); [Bae et al. 2017](#)), followed by Paleolithic migration to East Asia ~30 ka, inferred from genetic data ([Yan et al. 2014](#); [Hallast et al. 2015](#)). There was also an expansion of Neolithic farmers and Bronze Age migrations from southern China to MSEA, which contributed to the present-day gene pool of modern MSEA people, for example,

Thais and Laotians ([Higham 2014, 2017](#); [Lipson et al. 2018](#); [McColl et al. 2018](#)). Additional migrations during the historical period from neighboring countries ([Penth 2000](#); [Schliesinger 2000](#)) have further enhanced ethnolinguistic diversity.

The census size for Thailand was ~68.41 million in 2017 and for Laos was ~6.76 million in 2016 ([Simons and Fennig 2018](#)). There are five linguistic families distributed in these two countries. Although the Tai-Kadai (TK) language is widely spread in southern China and MSEA, it is concentrated in present-day Thailand and Laos as it is a major language spoken by Thais (90.5%) and Laotians (67.7%). Austroasiatic (AA) speakers are next most frequent, accounting for 4.0% in Thailand and 24.4% in Laos. In addition, this area is also inhabited by historical migrants who speak Sino-Tibetan (ST), Hmong-Mien (HM), and Austronesian languages (frequencies of 3.2%, 0.3%, and 2%, respectively, in Thailand; 3.1%, 4.8%, and 0% in Laos) ([Simons and Fennig 2018](#)).

© The Author(s) 2019. Published by Oxford University Press on behalf of the Society for Molecular Biology and Evolution.

This is an Open Access article distributed under the terms of the Creative Commons Attribution License (<http://creativecommons.org/licenses/by/4.0/>), which permits unrestricted reuse, distribution, and reproduction in any medium, provided the original work is properly cited.

Open Access

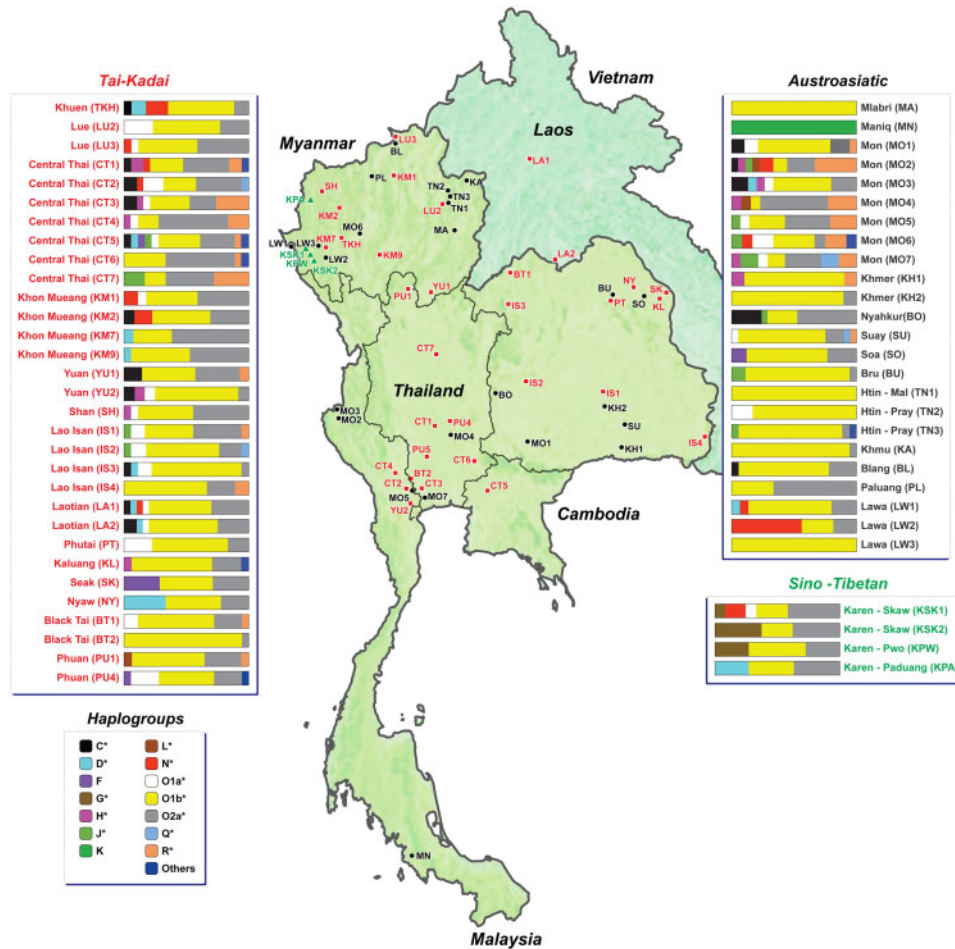


Fig. 1. Map showing sample locations and haplogroup distributions.

It is generally thought that AA languages were brought to the Thai/Lao region by Neolithic farmers from southern China, whereas TK languages were brought by a later, Bronze Age migration, also from southern China (Bellwood 2018). The Neolithic expansion was  $\sim 2-3$  ka before the expansion of TK languages; thus, the AA people were thought to be present before the TK expansion. The TK migration during the Bronze Age could have occurred via either demic diffusion (an expansion of TK people that brought both their genes and their language) or cultural diffusion (a language spread with minor movement of people). A genetic study on the origin of TK people supports a southern Chinese origin (Sun et al. 2013), whereas our previous studies of mitochondrial DNA (mtDNA) genome sequences support demic diffusion as the best explanation for the origin of the present-day Thai/Lao TK groups, although there is a strong signal of admixture between TK and AA groups in central Thailand (Kutanan et al. 2017; Kutanan, Kampuansai, Brunelli, et al. 2018). Although there is extensive ethnolinguistic diversity in the region, Thai/Lao populations can be generally categorized based on geography as either hill tribes or lowlanders. Nine ethnic groups, consisting of  $\sim 700,000$  people, are officially identified as hill tribes in Thailand: the AA-speaking Lawa, Htin, and Khmu; the HM-speaking Hmong and luMien; and the ST-speaking Karen, Lahu, Akha, and Lisu.

The Akha, Lisu, Hmong, luMien, Lawa, and Khmu are strongly patrilocal (i.e., the wife moves to the residence of her husband after marriage), whereas the Lahu, Karen, and Htin are strongly matrilocal. The lowlanders are neither strongly patrilocal nor matrilocal (Schliesinger 2000, 2001; Penth and Forbes 2004).

Previous studies have reported an influence of postmarital residence pattern on genetic variation in northern Thai hill tribes, with lower within-population genetic diversity coupled with greater genetic heterogeneity among populations for patrilocal groups than for matrilocal groups for the male-specific portions of the Y chromosome (MSY), whereas the opposite pattern is observed for mtDNA (Oota et al. 2001; Besaggio et al. 2007). However, these previous studies compared genetic variation between partial mtDNA sequences (hypervariable regions of the control region) and Y chromosomal short tandem repeats (Y-STRs); it would be informative to investigate more complete genetic data from these groups.

The MSY are paternally inherited and exhibit lineages specific to populations/geographic regions, making the MSY an informative tool for reconstructing paternal genetic history and demographic change (Yan et al. 2014; Barbieri et al. 2016). However, to date, there have been few MSY studies of MSEA and almost all of them employed Y-STRs (Cai et al. 2011; Kutanan et al. 2011; Brunelli et al. 2017) and also defined

haplogroups by genotyping assays, which are thus biased in terms of the haplogroups detected, and cannot uncover new sublineages. Analyzing partial sequences of the MSY and complete mtDNA genome sequences provides more insight into genetic history, especially sex-biased practices that can influence genetic variation, as well as the role of geography and language (Arias et al. 2018; Bajic et al. 2018; Kutanan, Kampuansai, Changmai, et al. 2018).

We have previously carried out comprehensive studies of the maternal genetic history of the Thai/Lao region, based on 1,823 complete mtDNA genome sequences (Kutanan et al. 2017; Kutanan, Kampuansai, Brunelli, et al. 2018; Kutanan, Kampuansai, Changmai, et al. 2018). In order to investigate the paternal genetic variation and demographic history, here, we investigate ~2.3 mB of MSY sequence in a subset of the above individuals, comprising 928 sequences from 59 populations. We compare and contrast the MSY and mtDNA results, with a focus on the patrilineal versus matrilineal hill tribes, the AA-speaking versus TK-speaking groups, and the various geographic regions (northern Thailand, central Thailand, and northeastern Thailand and Laos). We also use demographic modeling to address the role of demic versus cultural diffusion versus admixture in the origins of the major TK groups in each Thai/Lao region and contrast the results based on the MSY to previous results based on mtDNA. Our MSY sequencing results provide new insights into the paternal genetic history of MSEA and indicated contrasting paternal and maternal histories in this region.

## Results

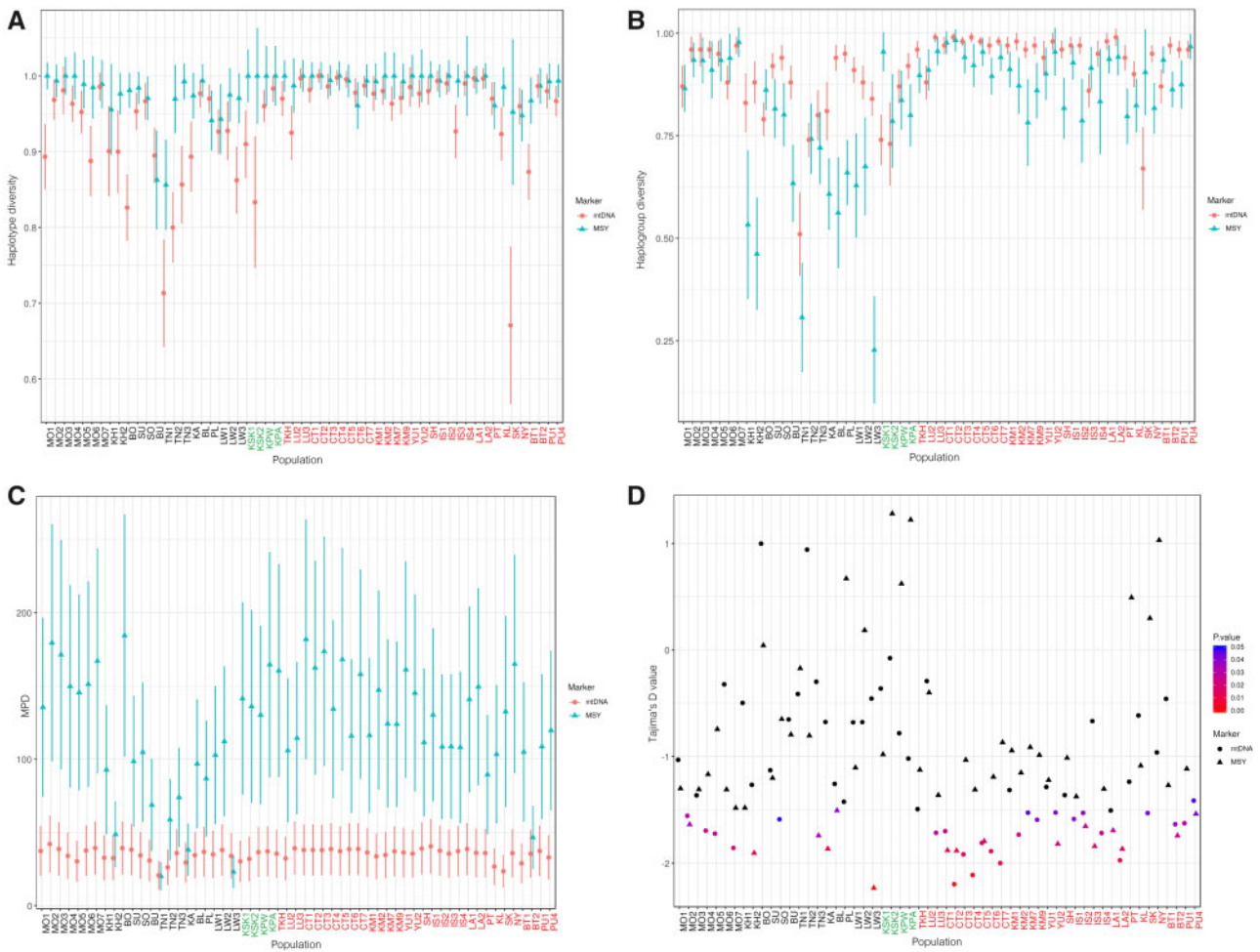
We generated 914 sequences of ~2.3 mB of the MSY, which combined with 14 published sequences brings the total to 928 MSY sequences belonging to 59 populations from Thailand and Laos (fig. 1 and supplementary table 1, Supplementary Material online). There are 816 haplotypes defined by 8,160 polymorphic sites, with mean coverages ranging from 4× to 109× (overall average coverage = 23×). Among the 928 MSY sequences, there are 92 specific haplogroups, belonging mostly to two main MSY lineages (O1b\* and O2a\*), that contribute substantially to the paternal genetic makeup of Thailand and Laos. There are several subclades of O1b\*; the most frequent (50.54%) is O1b1a1a\* or O-M95\*, which occurs in almost half of the AA groups with a very high frequency (>70%), that is, KH1-KH2, KA, BU, BL, SU, TN1-TN3, MA, and LW3 (fig. 1 and supplementary table 2, Supplementary Material online). The correspondence analysis (based on haplogroup frequencies) also supports the divergence of these AA-speaking groups in agreement with the other results mentioned later, with many O1b\* sublineages, for example, O1b1a1a1b1a (O-B426) and O1b1a1a1a1a (O-F2758) (supplementary fig. 1, Supplementary Material online). O2a\* or O-M324\* is the second most frequent haplogroup (25.86%) and has a relatively high frequency (>40%) in some AA and TK groups, and all ST-speaking Karen. Additional minor non-SEA-specific haplogroups were also observed, for example, haplogroup N\*, found in the Lawa groups, and haplogroups R\*, H\*, and J\*, which support associations

between India and the Mon, and genetic connections between Mon and TK groups (fig. 1 and supplementary fig. 1, Supplementary Material online). Further details on haplogroup distribution are provided in supplementary table 2 and text, Supplementary Material online.

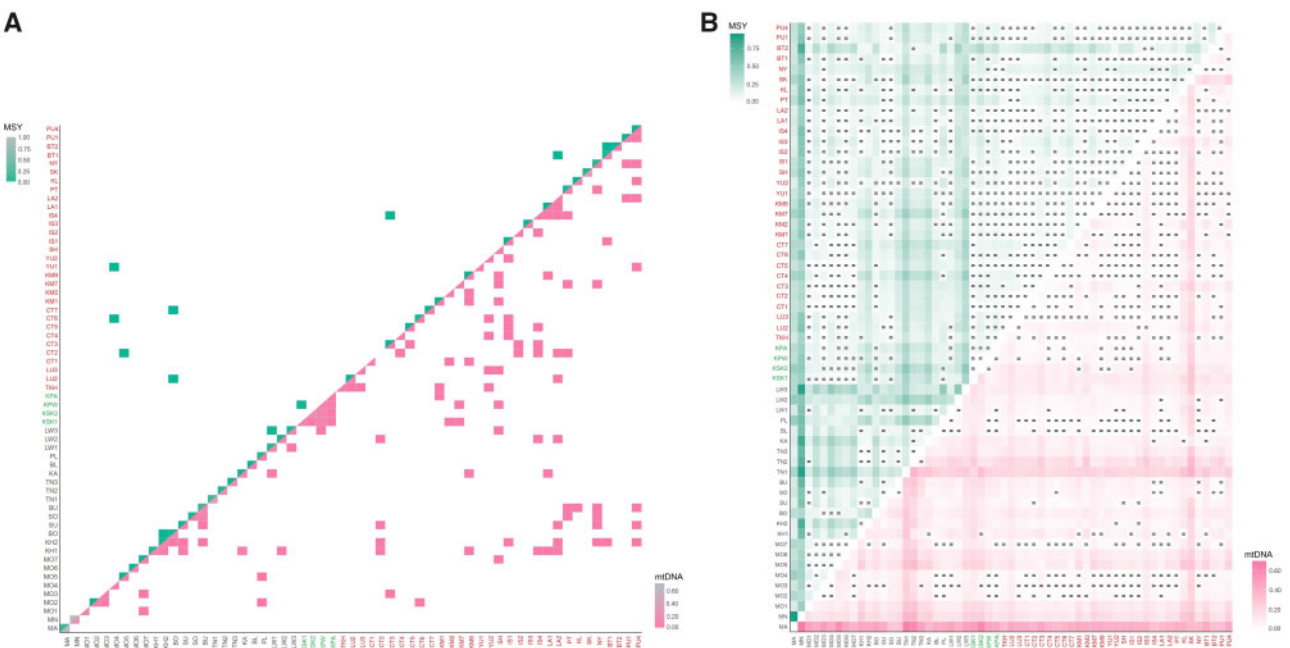
## Genetic Diversity and Structure

Generally, the AA populations show lower genetic diversity values than the TK and ST groups for the MSY, in agreement with the mtDNA results (fig. 2A–C) (Mann–Whitney *U* tests between AA and TK for MSY: *h*:  $Z = 3.37$ ,  $P < 0.01$ ; mean number of pairwise difference [MPD]:  $Z = 2.40$ ,  $P < 0.05$ ; haplogroup diversity:  $Z = 3.74$ ,  $P < 0.01$  and for mtDNA: *h*:  $Z = 4.33$ ,  $P < 0.01$ ; MPD:  $Z = 1.47$ ,  $P > 0.05$ ; haplogroup diversity:  $Z = 4.37$ ,  $P < 0.01$ ). After the Maniq (MN), who have no MSY variation, and the Mlabri (MA), who have no mtDNA variation, the Htin (TN1), Lawa (LW3), and Bru (BU) show very low diversity values of MSY, whereas the Htin (TN1–TN3), Khmer (KH2), and Seak (SK) show low mtDNA diversity (fig. 2A–C). In contrast to the other AA groups, the Mon (MO1–MO7) show higher levels of both MSY and mtDNA diversity than other AA groups (Mann–Whitney *U* tests between AA and Mon for MSY: *h*:  $Z = -3.33$ ,  $P < 0.01$ ; MPD:  $Z = -3.30$ ,  $P < 0.01$ ; haplogroup diversity:  $Z = -3.75$ ,  $P < 0.01$  and for mtDNA: *h*:  $Z = -1.94$ ,  $P > 0.05$ ; MPD:  $Z = -2.03$ ,  $P < 0.05$ ; haplogroup diversity:  $Z = -2.79$ ,  $P < 0.01$ ). LW3 showed very low MSY haplogroup diversity (fig. 2B) and MPD values (fig. 2C), and a significantly low Tajima's *D* value (fig. 2D), suggesting recent paternal expansion in this group, but the converse trend (rather high diversity) for mtDNA. Interestingly, a significantly negative Tajima's *D* value was observed more frequently in the TK than the AA groups for both the MSY and mtDNA (MSY,  $P < 0.05$ : 10/31 for TK vs. 6/24 for AA; mtDNA,  $P < 0.05$ : 20/31 for TK vs. 5/24 for AA) (fig. 2D), suggesting a stronger signal of recent population expansion in TK groups; no significant Tajima's *D* values were observed in any of the ST-speaking Karen groups. The Nyahkur (BO), who speak a Mon language, show the highest MPD value for the MSY (fig. 2C), which might indicate paternal gene flow with other populations; this is supported by the BO having the highest number of shared MSY haplotypes (three haplotypes) with other populations (fig. 3A). MO3 and MO4 have shared MSY haplotypes with the TK-speaking groups (CT2, CT6, and YU1), reflecting their genetic connection. In the mtDNA, apart from the AA-speaking Palaung (PL), the Mon (MO2, MO3, and MO7) also share haplotypes with the central Thai (CT3 and CT6) and Shan (SH) (fig. 3A).

The Analysis of Molecular Variance (AMOVA) indicates that the variation among populations (within group) accounts for 11.12% of the total MSY genetic variance (table 1). There is greater genetic heterogeneity within the AA group (20.01%,  $P < 0.01$  and 18.49%,  $P < 0.01$  without MN, the hunter–gatherer group from southern Thailand) than among the TK (4.48%,  $P < 0.01$ ) and ST-speaking Karen groups (2.29%,  $P > 0.01$ ). For the AA group with more than one population sampled, the greatest within-group variation by far was among the three Lawa populations



**FIG. 2.** Genetic diversity values of MSY and mtDNA in the studied populations, excluding the Maniq (MN) and Mlabri (MA): haplotype diversity (A), haplogroup diversity (B), MPD (C), and Tajima's *D* values (D). More information and all genetic diversity values are provided in [supplementary table 1, Supplementary Material](#) online.



**FIG. 3.** Relative shared haplotypes (A) and heat plot of  $\Phi_{st}$  (B) between studied populations for the MSY and for mtDNA.

**Table 1.** AMOVA Results.

Groups	Number of Groups	Number of Populations	Percent Variation					
			Within Populations		Within Groups		Among Groups	
			MSY	mtDNA	MSY	mtDNA	MSY	mtDNA
Total	1	59 (58)	88.88 (89.46)	91.51	11.12* (10.54*)	8.55*		
Language	3	59 (58)	88.21* (98.05*)	91.20*	10.16* (1.96*)	8.18*	1.63* (−0.01)	0.62*
Austroasiatic	1	24 (23)	79.99 (81.51)	85.97	20.01* (18.49*)	14.03*		
Mon	1	7	96.08	93.10	3.92*	6.90*		
Htin	1	3	88.47	74.29	11.53*	25.71*		
Lawa	1	3	65.57	92.22	34.43*	7.78*		
Sino-Tibetan (Karen)	1	4	97.71	93.49	2.29	6.51*		
Tai-Kadai	1	31	95.52	95.67	4.48*	4.33*		
Central Thai	1	7	98.53	98.36	1.47	1.64*		
Khon Mueang	1	4	101.83	95.80	−1.83	4.20*		
Lao Isan	1	4	98.16	97.69	1.84	2.31*		
Geography	6 (5)	59 (58)	88.27* (98.07*)	91.40*	9.35* (2.02*)	8.40*	2.38* (−0.09)	0.20*
Northern	1	26	85.51	88.84	14.49*	11.16*		
Northeastern	1	16	96	91.29	8.00*	8.71*		
Central	1	11	94.61	95.86	5.39*	4.14*		
Western	1	3	93.97	99.11	6.03*	0.89		

NOTE.—The numbers in parentheses show the percent variation of MSY by excluding the Maniq (MN) and asterisks indicate significant level ( $P < 0.01$ ).

(34.43%,  $P < 0.01$ ), whereas the seven Mon populations showed very low (albeit still significant) within-group variation (3.92%,  $P < 0.01$ ) (supplementary fig. 2, Supplementary Material online). Very low within-group variation was also observed for the central Thai groups from central Thailand (1.47%  $P > 0.01$ ), Khon Mueang groups from northern Thailand (−1.83%,  $P > 0.01$ ), and Lao Isan groups from northeastern Thailand (1.84%,  $P > 0.01$ ), indicating overall genetic homogeneity among these major TK-speaking groups. In agreement with the MSY, larger mtDNA variation is observed in the AA groups (14.03%,  $P < 0.01$ ) than the ST (6.51%,  $P < 0.01$ ) and TK groups (4.33%,  $P < 0.01$ ), but interestingly the largest within-group variation is not among the Lawa (7.78%,  $P < 0.01$ ) but rather among the Htin populations (25.71%,  $P < 0.01$ ). In contrast to the MSY, each of the TK groups with more than one population sampled showed significant within-group differences for mtDNA, especially the Khon Mueang (4.20%,  $P < 0.01$ ) (supplementary fig. 2, Supplementary Material online). In sum, we observed different patterns of MSY versus mtDNA for the different language groups. The among-population variation within linguistic groups is larger for the MSY (20.01%,  $P < 0.01$ ) than for mtDNA (14.03%,  $P < 0.01$ ) for AA groups, but about the same for TK groups (4.48%,  $P < 0.01$  for MSY and 4.33%,  $P < 0.01$  for mtDNA), and the ST groups have larger among-population variation for mtDNA (6.51%,  $P < 0.01$ ) than for the MSY (2.29%,  $P < 0.01$ ) (table 1 and supplementary fig. 2, Supplementary Material online). Thus, there are different patterns of MSY versus mtDNA differentiation for these three language families.

Although there is more variation among groups defined by geographic location (2.38%,  $P < 0.01$ ) than by language family (1.63%,  $P < 0.01$ ) (table 1), there is much more MSY variation among populations within the same group than among groups defined either by geographic or by linguistic criteria. Moreover, when the divergent MN population of hunter-

gatherers from southern Thailand is removed from the analysis, then the among-group component is no longer significant for either geographic location or language family (−0.09%,  $P > 0.01$  for geography; −0.01%,  $P > 0.01$  for language), and the total variation among populations within group reduces to 10.54%. Thus, neither geography nor language family is a good predictor of the MSY genetic structure of Thai/Lao populations, indicating that these two factors are not important in the broad view (table 1).

There are significant correlations between matrices of MSY genetic and geographic distance, estimated by Mantel tests, for all three types of geographic distances, that is, great circle distance ( $r = 0.3381$ ,  $P < 0.01$ ), resistance distance ( $r = 0.5418$ ,  $P < 0.01$ ) and least-cost path distance ( $r = 0.3912$ ,  $P < 0.01$ ). However, the correlations are no longer significant when the MN group is removed from the analysis: great circle distance ( $r = 0.0125$ ,  $P > 0.05$ ), resistance distance ( $r = -0.0446$ ,  $P > 0.05$ ) and least-cost path distance ( $r = 0.0139$ ,  $P > 0.05$ ). In contrary, no significance was detected ( $P > 0.05$ ) between matrices of mtDNA genetic distance and geographic distances with and without MN (great circle distance:  $r = 0.0776$  and  $r = -0.0323$ ), resistance distance ( $r = 0.1433$  and  $r = -0.1105$ ), and least-cost path distance ( $r = 0.0997$  and  $r = -0.0253$ ).

To identify and describe population clustering based on multivariate analysis, discriminant analysis of principal components (DAPC) was carried out. This analysis attempts to maximize among-groups genetic differentiation and minimize within-group genetic variation; the results showed considerable overlap among groups defined by either language family or geographic location in both MSY and mtDNA (supplementary fig. 3, Supplementary Material online). In addition, the groupings by population and ethnicity of MSY data revealed the largest discrimination to be among some AA-speaking groups, that is, all Lawa groups (LW1–LW3), Htin (TN1), and Blang (BL), whereas all Htin groups (TN1, TN2,

and TN3), Mlabri (MA), TK-speaking Seak (SK), and ST-speaking Karen (KSK1, KSK2, and KPW) are differentiated from the others for mtDNA, emphasizing contrasting genetic pattern between MSY and mtDNA for Htin, Mlabri, Lawa, Blang, Seak, and Karen.

In sum, all results indicate lower genetic diversity of the AA groups than the TK and ST groups, except the Mon and Nyahkur, who exhibit high genetic diversity. The AA groups also show greater genetic heterogeneity than the TK and ST groups.

### Postmarital Residence and Genetic Diversity

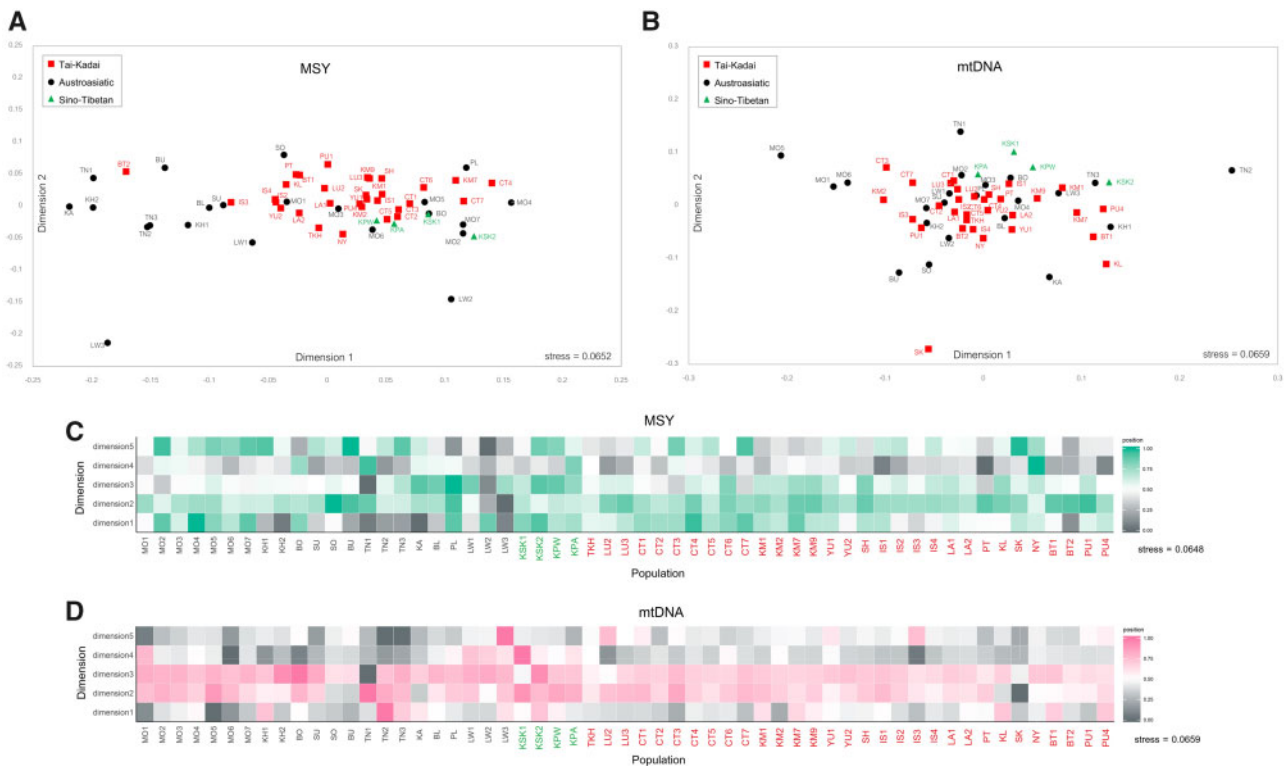
We studied five highlander groups: four hill tribes (Karen, Htin, Lawa, and Khmu) and the Palaung, another minority group in the mountainous area of northern Thailand but not officially recognized as a hill tribe. The Khmu (KA), Lawa (LW1, LW2, and LW3), and Palaung (PL) groups practice patrilocality, whereas the Htin (TN1, TN2, and TN3) are matrilocal, as are the ST-speaking Karen (KSK1, KSK2, KPA, and KPW). If residence pattern is influencing genetic variation, then lower within-population genetic diversity coupled with greater genetic heterogeneity among populations is expected for patrilocal groups than for matrilocal groups for the MSY, whereas the opposite pattern is expected for mtDNA (Oota et al. 2001). The MSY  $h$  and MPD values are higher for matrilocal groups, but not significantly (Mann–Whitney  $U$  tests:  $h$ :  $Z = 1.4616$ ,  $P > 0.05$ ; MPD:  $Z = 0.9744$ ,  $P > 0.05$ ); however, haplogroup diversity is significantly higher for the matrilocal groups (Mann–Whitney  $U$  tests:  $Z = 2.1112$ ,  $P < 0.05$ ) (supplementary fig. 4, Supplementary Material online). For mtDNA, genetic diversity values are higher for patrilocal than for matrilocal groups, but the differences are not statistically significant (Mann–Whitney  $U$  tests:  $h$ :  $Z = -0.9744$ ,  $P > 0.05$ ; MPD:  $Z = -0.8120$ ,  $P > 0.05$ ; haplogroup diversity:  $Z = -1.864$ ,  $P > 0.05$ ) (supplementary fig. 4, Supplementary Material online). Notably, TN1 and LW3 exhibit very low within-population diversity for the MSY, for example, MPD = 20.07 and 23.07, compared with the average MPD (121.11), whereas TN1 and TN2 (20.69 and 26.14) show lower MPD than average (35.09) for mtDNA (supplementary table 1, Supplementary Material online). For genetic differences between-populations, the patrilocal Khmu, Lawa, and Palaung have significantly higher genetic differentiation for the MSY than for mtDNA (average  $\Phi_{st} = 0.3109$  for MSY and 0.0774 for mtDNA) (Mann–Whitney  $U$  tests:  $Z = 3.5907$ ,  $P < 0.01$ ), whereas the matrilocal groups (Htin and Karen) also show higher average  $\Phi_{st}$  for MSY (0.1859) than for mtDNA (0.1553), but these are not significantly different (Mann–Whitney  $U$  tests:  $Z = 0.3270$ ,  $P > 0.05$ ). Contrasting genetic differences for the MSY versus mtDNA of Lawa, Htin, and Karen are clearly seen in the multidimensional scaling (MDS) and DAPC plots (fig. 4A and B and supplementary fig. 3, Supplementary Material online). Much stronger contrasting between-group variation is seen in the AMOVA results (Lawa: 34.43% for MSY and 7.78% for mtDNA; Htin: 11.53% for MSY and 25.71% for mtDNA; Karen: 2.29% for MSY and 6.51% for mtDNA) (table 1 and supplementary fig. 2, Supplementary Material online).

However, in general, the AA-speaking groups, whether identified as hill tribes or as other minorities, are patrilocal groups. The AMOVA result indicates that the variation among AA populations is higher in MSY (20.01%) than mtDNA (14.03%), in accordance with expectations if residence pattern is influencing genetic variation. Conversely, the TK populations, where neither patrilocal nor matrilocal residence is preferred, exhibit similar among-population variances for the MSY (4.48%) and mtDNA (4.33%) (table 1 and supplementary fig. 2, Supplementary Material online). Overall, there does seem to be some impact of postmarital residence on the patterns of genetic diversity.

### Genetic Relatedness among Populations

The genetic distance and MDS analyses based on MSY and mtDNA indicate that the MN and MA are highly diverged from the other populations for the MSY and mtDNA, respectively (supplementary fig. 5, Supplementary Material online). The MA and MN also show large differences from the other populations in the heat plots of  $\Phi_{st}$  values (fig. 3B). However, in general, both MSY and mtDNA results show relatively larger genetic heterogeneity of the AA groups versus genetic homogeneity of the TK and ST groups (fig. 3B). The Mantel test of  $\Phi_{st}$  values showed a significant correlation between the MSY and mtDNA  $\Phi_{st}$  matrices ( $r = 0.4506$ ,  $P < 0.01$ ). After excluding these MA and MN as outliers, the MDS for the MSY showed that almost all AA-speaking groups are located along the edges of the plot, whereas most of the TK groups cluster in the center of the plot (fig. 4A), further supporting genetic heterogeneity of the AA and homogeneity of the TK populations. Interestingly, the SEA-specific O-M95\* and O-M324\* haplogroups (with several sublineages) differentiate the studied populations into at least two main paternal sources, and the frequencies of these two haplogroups correspond to the major differentiation in the MDS plot (fig. 4A). O-M95\* is at high frequency in the populations on the left of the plot and gradually decreases to very low frequency in the populations on the right side in the first dimension, whereas the O-M324\* frequency runs opposite to the O-M95\* cline: O-M324\* is at higher frequency in populations located on the right of the plot and decreases in frequency toward the left side (fig. 4A). The MDS plot and heat plot of MSY also indicates some Mon groups (MO1, MO3, MO5, and MO6) are close to the cluster of TK groups in the center of the plot (fig. 4A and C), indicating a close genetic relationship. In addition, non-SEA haplogroups lineages, for example, R\*, H\*, and J\*, provide more support for genetic connections between Mon and Central Thais.

For the MDS based on mtDNA (fig. 4B), the Mon generally showed genetic affinities with the TK groups in the center of the plot, with the exception of MO1, MO5, and MO6, which differ from the other Mon groups, as can be also seen in the MDS plot and heat plot (fig. 4B and D). Overall, we observe more genetic heterogeneity of the AA groups than the other



**Fig. 4.** The two-dimensional MDS plot and five-dimensional MDS heat plot based on the  $\Phi_{st}$  distance matrix for 57 populations (after removal of Maniq and Mlabri) of MSY (A and C) and mtDNA (B and D).

linguistic groups and there are contrasting patterns of genetic relationships for the MSY versus mtDNA.

### Genetic Relatedness between Thai/Lao and Other Asian Populations

The MDS based on the MSY  $\Phi_{st}$  matrix of 73 populations from across Asia revealed that, in general, population clustering largely reflects linguistic affiliation (fig. 5), with some exceptions. In the first and second dimension, the AA populations are the most diversified, with the PL and MN appearing as outliers. There is one cluster of AA populations on the left, which also includes one TK group (BT2); the other AA populations are scattered along the main axis of the plot. Some Mon groups (MO2, MO4, and MO7) are relatively close to Indian and ISEA populations, indicating potential connections. Two central Thai groups (CT4 and CT7) are also relatively close to the Indian populations. The ST populations (Karen, Han Chinese, and Burmese) are rather close. The ISEA and Papuan populations are in closer proximity to South Asian populations (Indian, Bengali, and Punjabi). Generally, the haplogroup profile indicates genetic affinities between the Mon and South/Central Asian groups, which is consistent with the MDS plots (fig. 5) and results from previous mtDNA haplogroup analyses (Kutanan et al. 2017; Kutanan, Kampuansai, Brunelli, et al. 2018).

### The Expansion of Male Lineages

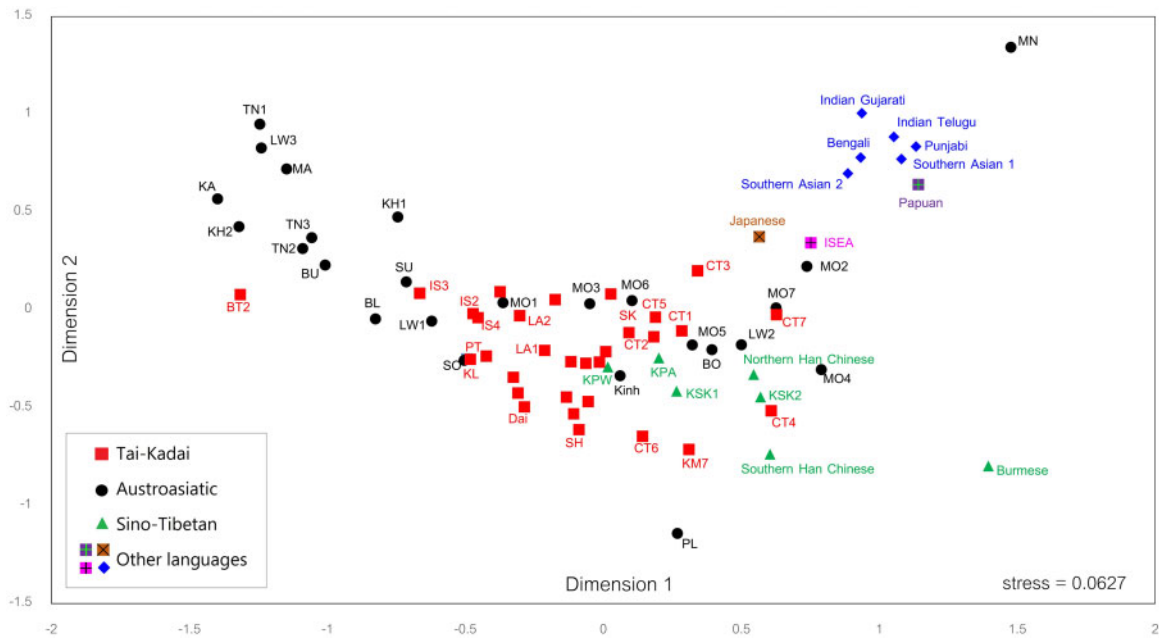
The Bayesian Skyline Plots (BSPs) of effective population size change ( $N_e$ ) over time in each group reveal overall five

different trends (fig. 6). The most common trend, found in Mon, Khmer, Htin, Central Thai, and Black Tai, showed  $N_e$  increasing gradually or remaining constant during 40–60 ka until a decline  $\sim$ 5–7 ka, followed by rapid growth  $\sim$ 5 ka and then a decrease  $\sim$ 2.0–2.5 ka. The other trends differ from the first trend as follows: no population reduction  $\sim$ 2.0–2.5 ka but population size either increases (Khon Mueang and Yuan) or remains stable (Lao Isan and Laotian); the Lue and Phuan show two increases in  $N_e$ , at about  $\sim$ 5 ka and  $\sim$ 10 ka; the Lawa show a stable population size since  $\sim$ 30 ka and then a decline during the last 2 ka with a sudden increase  $\sim$ 1 ka; and the Karen differ only slightly from the common trend with a population increase  $\sim$ 1 ka.

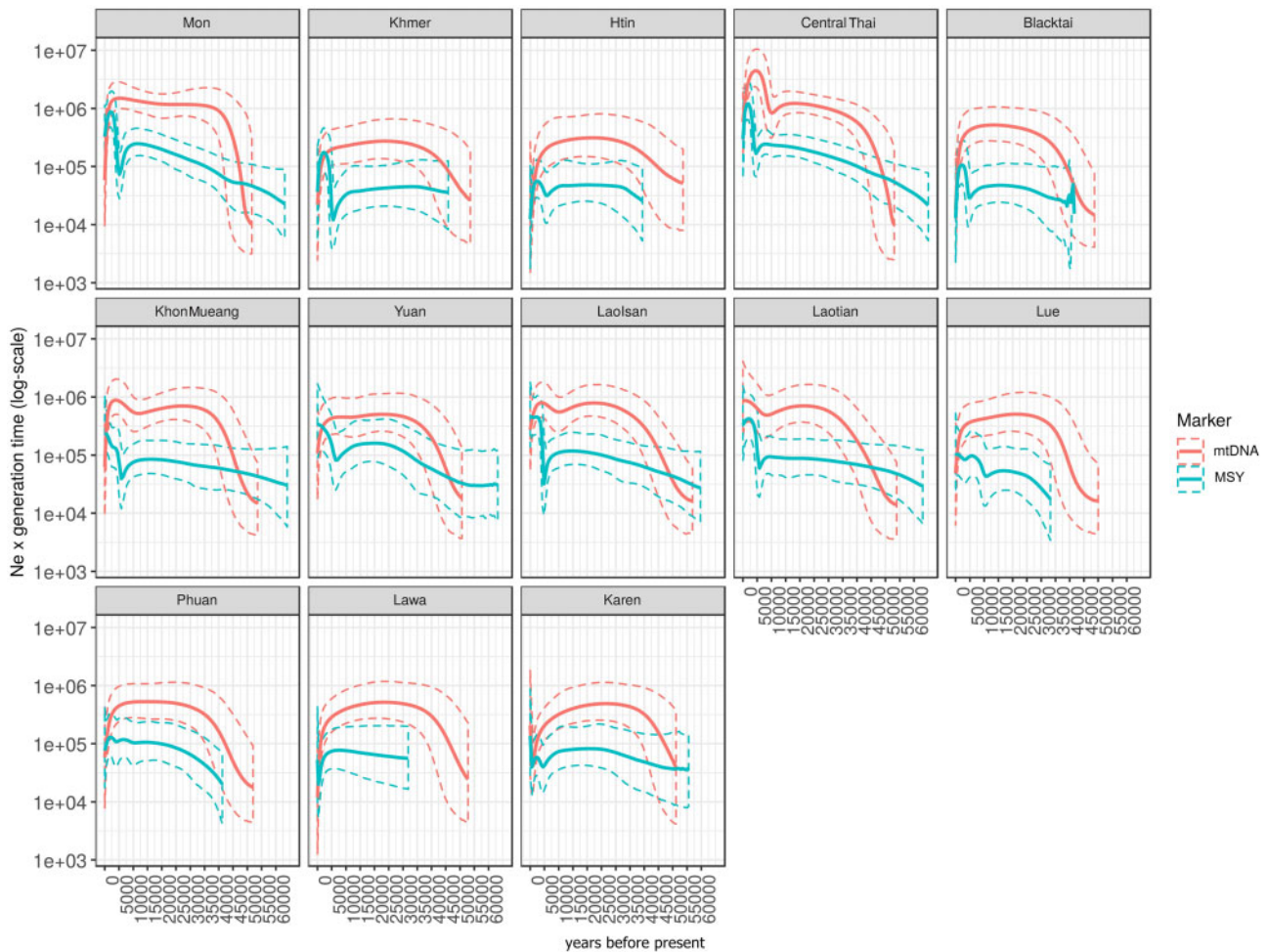
By contrast, the BSP based on mtDNA sequences for each ethnicity show three common trends (fig. 6). The first trend is an increase in  $N_e$  during 40–50 ka, followed by stability and then decrease  $\sim$ 2 ka, which was observed in Mon, Htin, Lawa, Khmer, Yuan, Phuan, and Lue. The second pattern, shown by the Khon Mueang is an increase in  $N_e$   $\sim$  40–50 ka, followed by stability and then increase again  $\sim$ 10 ka, followed by a decline  $\sim$ 2 ka. The Central Thai, Lao Isan, and Laotian show the third trend, in which population increases occur  $\sim$ 40–50 and  $\sim$ 10 ka. In general, the BSP by ethnicity indicated lower effective population sizes for the MSY than for mtDNA (fig. 6).

We also plotted the BSP of several Asian populations from published MSY data (Karmin et al. 2015; Poznik et al. 2016) (fig. 7). Almost all of the MSEA and East Asian populations, that is, Kinh, Northern Han, Southern Han, and Japanese

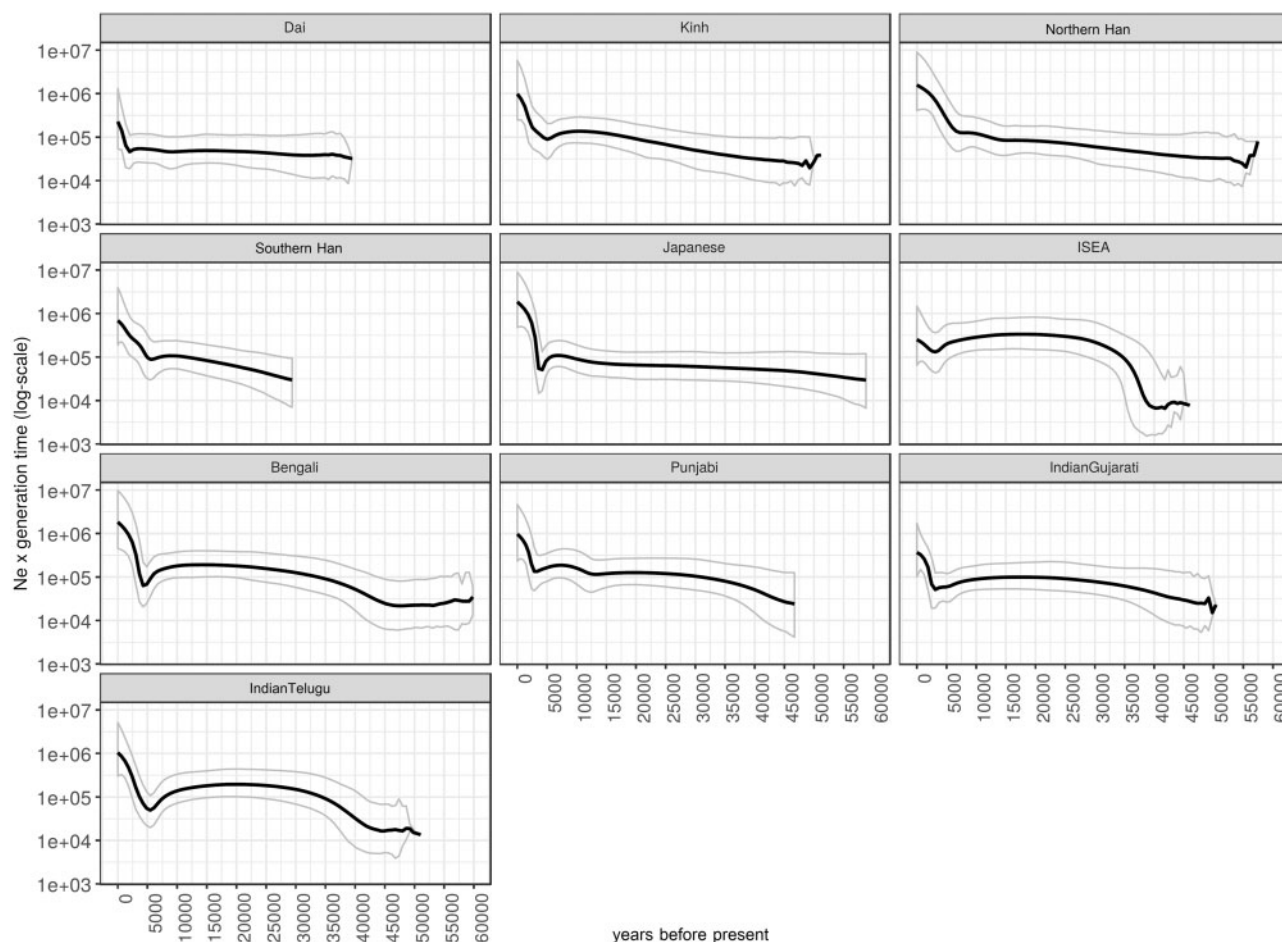




**Fig. 5.** The two-dimensional MDS plot based on the  $MSY \Phi_{st}$  distance matrix for 73 populations. Population details are listed in [figure 1](#) and [supplementary tables 1 and 7, Supplementary Material](#) online.



**Fig. 6.** The BSPs based on the MSY and mtDNA of 13 ethnicities from Thailand and Laos; Mon, Khmer, Htin, Central Thai, Black Tai, Khon Mueang, Yuan, Lao Isan, Laotian, Lue, Phuan, Lawa, and Karen. Solid lines are the median estimated effective population size ( $y$  axis) through time from the present in years ( $x$  axis). The 95% highest posterior density limits are indicated by dotted lines.



**FIG. 7.** The BSPs of Asian populations. Solid lines are the median estimated paternal effective population size (y axis) through time from the present in years (x axis). The 95% highest posterior density limits are indicated by dotted lines.

show a pronounced increase of the MSY  $N_e$  during  $\sim 4$ – $6$  ka, except the Xishuangbanna Dai, in which there is an increase  $\sim 2$  ka. Around 5 ka, the Japanese show a decrease in  $N_e$  before a sudden increase, suggesting a bottleneck prior to demographic expansion. Interestingly, the ISEA population shows a large increase in  $N_e \sim 35$ – $40$  ka and a smaller increase  $\sim 2.5$ – $3$  ka. The South Asian populations, that is, Bengali, Punjabi, and Indians, also show two pulses of population increase at about the same times. The Punjabi also show an additional small increase in  $N_e$  change during  $\sim 12$  ka.

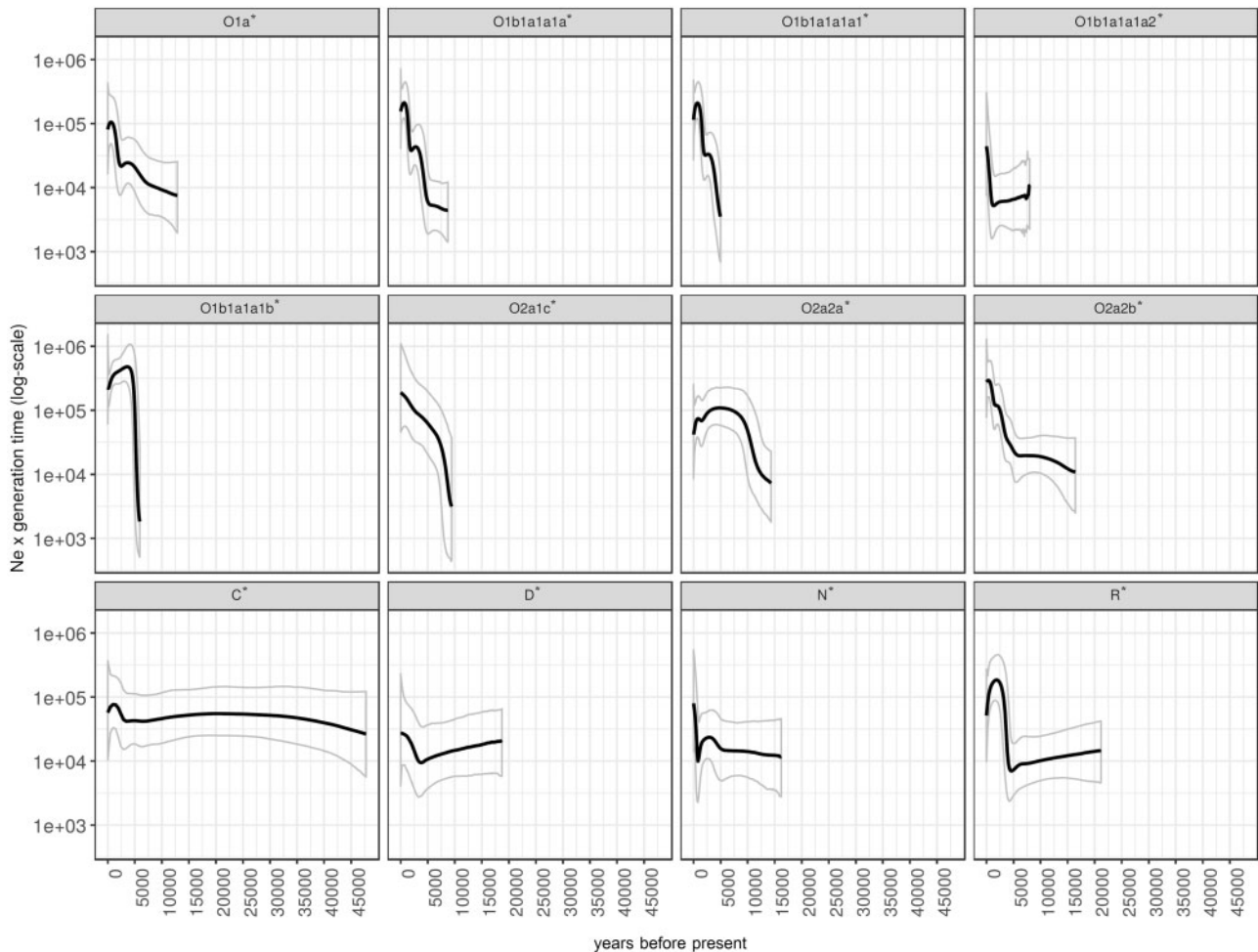
The BSP by each major MSY haplogroup show four pulses of paternal  $N_e$  increases, at  $\sim 9$ – $11$  ka,  $\sim 5$  ka,  $\sim 2.0$ – $2.5$  ka, and  $\sim 1.0$  ka (fig. 8), in agreement with the plot by ethnicity. The early Holocene  $N_e$  increment is obviously noticed in O2a1c\* and O2a2a\*, whereas the  $N_e$  growth  $\sim 5$  ka is observed in O1b1a1a1b\* and R\*. Haplogroup O1a\*, C\* and D\* show expansions in  $N_e \sim 2.0$ – $2.5$  ka and haplogroup N\* shows a recent expansion  $\sim 1.0$  ka. In addition, there are two expansion times for O1b1a1a1a\* and O2a2b\* ( $\sim 5$  and  $\sim 2$  ka).

### Demographic Models

Previously, we used mtDNA genome sequences and demographic modeling to test different hypotheses about the origins of TK groups. Specifically, we tested whether different TK groups were primarily related to local AA groups (reflecting

cultural diffusion, i.e., an AA group switching to a TK language), to a TK group from southern China (reflecting demic diffusion, i.e., spread of TK languages via migration from southern China), or were related to both (reflecting admixture between an incoming TK group from southern China and a local AA group). We found that the Khon Mueang (from northern Thailand), Lao Isan (from northeastern Thailand), and Laotian most likely originated via demic diffusion from southern China without substantial gene flow from AA groups (Kutanan et al. 2017). However, for the central Thai, the most likely scenario was demic diffusion with a very low level of gene flow between central Thai and Mon groups (Kutanan, Kampuansai, Brunelli, et al. 2018). Here, we use the same approach to test three demographic scenarios concerning the paternal origins of these major Thai groups (supplementary fig. 6, Supplementary Material online).

For the Khon Mueang (KM) people (Test 1), the highest posterior probability (0.80) and rather highly selected classification trees (0.58) were found for the demic diffusion model (supplementary table 3, Supplementary Material online). By contrast, the cultural diffusion model is the most likely scenario for the Lao and central Thai groups. Both the combined Laotian (LA) and Lao Isan (IS) data sets (Test 2) and the separate LA data set (Test 3) weakly support the cultural diffusion model (for Test 2: posterior probability = 0.56



**Fig. 8.** The BSPs for each major haplogroup. Solid lines are the median estimated paternal effective population size ( $y$  axis) through time from the present in years ( $x$  axis). The 95% highest posterior density limits are indicated by dotted lines.

and selected classification tree = 0.37 and for Test 3: posterior probability = 0.56 and selected classification tree = 0.39). The IS data set (Test 4) supports cultural diffusion (with the present-day IS groups descended from local Khmer [KH] with the highest posterior probability [0.71] and classification trees selected slightly more often than for the other models [0.49]). For Test 5 (the central Thai [CT] data set), the cultural diffusion model had the highest posterior probability at 0.58 and was selected slightly more often among the classification trees (0.50) than the other models. However, a Principal Component Analysis plot shows that based on the first two PCs the observed data fall within the distributions simulated under the three models in only Test 4, whereas the other data sets fall within the simulated distributions for PCs 3 and 4, suggesting that there is low efficiency to reconstruct the variability of the observed data (supplementary fig. 7, Supplementary Material online). The parameter estimation for the best performing models in all five tests was able to obtain point estimates for each of the simulated effective population sizes (supplementary table 4, Supplementary Material online). However, the posterior distributions were generally flat (supplementary fig. 8, Supplementary Material online). We also calculated the MSY  $\Phi_{st}$  and corrected

pairwise difference among groups of populations used in ABC tests to estimate their genetic relationships (supplementary table 5, Supplementary Material online). The KM are closer to the Dai than the local AA group (Test 1), the ethnic Lao and Laotian showed similar genetic differences to both Dai and AA groups (Test 2 and Test 3), whereas the CT groups (Test 5) have closer genetic relationships to the local AA group than to Dai. In contrast, mtDNA  $\Phi_{st}$  and corrected pairwise difference revealed that the KM and ethnic Lao are closer to the Dai than local AA, whereas the CT exhibited somewhat similar genetic distances to both Dai and AA. Overall, the simulations based on MSY sequences, compared with previous mtDNA simulation together with tests of genetic difference by  $\Phi_{st}$  and corrected pairwise differences, suggest different demographic histories for males and females in the region.

## Discussion

In order to gain more insights into MSEA genetic history, we here investigate the paternal genetic variation and structure by sequencing  $\sim 2.3$  mB of the MSY from representative ethnolinguistic groups from Thailand and Laos. In sum, most of the studied populations exhibit two major MSY

haplogroups, O-M324\* and O-M95\* in different proportions, indicating two major paternal sources. O-M324\* was widely spread in the TK groups, whereas O-M95\* is predominant in the AA groups. However, some TK populations (BT2 and IS3) and some AA populations (PL, BO and MO4) exhibited the opposite pattern (fig. 1 and supplementary table 2, Supplementary Material online). We also compared patterns of MSY variation with mtDNA in the same set of populations and found some similar results, for example, overall lower genetic diversity and greater heterogeneity of AA groups than of TK and ST groups, large differences between the Mon and the other AA groups, and genetic connections between the Mon and central Thai (figs. 2–4). However, in many respects, the patterns of MSY and mtDNA variation are different, suggesting contrasting paternal and maternal genetic histories. Here, we focus on three groups of populations with different cultural practices and histories that also stand out in the genetic analyses: the hill tribes, the AA-speaking Mon, and the major TK-speaking groups.

### Factors Influencing Contrasting Genetic Variation in the Hill Tribes

The hill tribes, who occupy the mountainous northern region of Thailand, are notable for their variation in patrilocal versus matrilocal residence pattern (Oota et al. 2001; Besaggio et al. 2007), as well as for their strong sense of group identity, which tends to isolate them from other groups (Schliesinger 2000; Nahhas 2007). If postmarital residence is influencing patterns of genetic variation, then the expectation is for larger between-group differences and smaller within-group diversity for patrilocal groups for the MSY, and the same trends for matrilocal groups for mtDNA. The first comparative study of mtDNA and MSY variation in patrilocal versus matrilocal groups was carried out in the northern Thai hill tribes and found a strong impact of postmarital residence on the mtDNA and MSY variation (Oota et al. 2001). However, previous studies compared genetic variation between partial mtDNA sequences and Y-STRs (Oota et al. 2001; Besaggio et al. 2007); here, we report the first comparison of mtDNA and MSY variation based on comparable sequence data.

Here, we analyzed the sequences of mtDNA genome and ~2.3 mB of the MSY of the Khmu, Palaung, and Lawa groups, who practice patrilocality, whereas the Htin are matrilocal, similar to the ST-speaking Karen. The within-population genetic diversity values are in agreement with expectations, that is, greater diversity of matrilocal than patrilocal groups for MSY and the opposite trend in mtDNA (supplementary fig. 4, Supplementary Material online). Moreover, genetic differentiation between populations also goes in the direction predicted by postmarital residence pattern. However, in many cases, the differences between patrilocal and matrilocal groups are not significant, indicating that other factors are also having an effect. In particular, the Htin (TN1) and Lawa (LW3) exhibit very low within-population diversity for the MSY, whereas the Htin (TN1 and TN2) also show lower diversity for the mtDNA (fig. 2A–C).

One factor in particular that could influence the within-population genetic diversity and between-population

differentiation is geographic isolation, which enhances genetic drift and inbreeding, thereby lowering within-population genetic diversity and increasing between-population differentiation. This could explain the very low internal diversity and high differentiation from other groups of some groups of Htin (TN1) and Lawa (fig. 4A and B and supplementary fig. 3, Supplementary Material online) that live in mountainous, isolated parts of northern Thailand. The Lawa furthermore favor intramarriage (Nahhas 2007) which would also reduce genetic variation in this group. The Htin (TN1) also show very low diversity and extreme divergence in genome-wide single nucleotide polymorphisms data (Xu et al. 2010) and both Htin (TN1–TN3) and Lawa (LW3) exhibit lower diversity and large differentiation in autosomal STRs (Kampuansai et al. 2017). Such drastic genetic drift effects could reduce the significance of the impact of postmarital residence on patterns of genetic diversity.

Moreover, these results are in keeping with previous observations that although the expected difference between patrilocal and matrilocal groups holds in some regions (Oota et al. 2001; Besaggio et al. 2007), in other regions patterns of mtDNA and MSY variation do not conform to expectations (Kumar et al. 2006; Arias et al. 2018). This is indeed to be expected given that many other factors, for example, other human cultures (e.g., linguistic exogamy), physical landscape, and subsistence strategies, influence patterns of genetic variation (Wilkins and Marlowe 2006; Chaix et al. 2007).

### Genetic Variation and Origin of the Mon

The Mon groups showed genetic differences from other AA populations but closer relatedness to the TK populations, especially the central Thai, in both MSY and mtDNA (figs. 2A, B, 3A, and B). Our previous simulation results, based on mtDNA, also supported admixture among the Mon and central Thai groups (Kutanan, Kampuansai, Brunelli, et al. 2018). In addition, some Mon groups (MSY: MO3, MO5, MO6 and mtDNA: MO2, MO3 and MO4) exhibit genetic affinities with the Karen (fig. 3B), reflecting genetic heterogeneity and contrasting genetic patterns between MSY and mtDNA. Admixture might be an important factor influencing the genetic structure of the lowland AA-speaking Mon. Archaeological evidence indicates that the Dvaravati civilization of the Mon was centered in present-day central Thailand and southern Myanmar and had expanded to a large part of MSEA during the sixth to seventh century AD (Diffloth 1984; Guillou 1999; Saraya 1999). After the intensification of Thai and Burmese kingdoms, the Mon in Myanmar were conquered by the Burmese during the 18th century AD; the ethnic Mon in Myanmar are currently concentrated in the Mon and Karen States (Pon Nya 2001). In Thailand, the present-day Mon are distributed in central Thailand and surrounding areas, with some groups living in the North and the Northeast. However, they are not considered to be the descendants of the ancient Mon Dvaravati civilization in Thailand, but rather political refugees that fled from Myanmar to Thailand during the 16th to 19th centuries AD (Ocharoen 1998). However, based on linguistic evidence, the remnants of the Dvaravati Mon population are now

considered a distinct ethnic group known as the Nyahkur (BO) whose communities are restrict found in hilly areas along the border between central and northeastern Thailand (Diffloth 1984). In contrary to linguistic evidence, the Nyahkur has no shared haplotype or related to any specific Mon groups, indicating their genetic differences. However, Nyahkur show genetic sharing in both MSY and mtDNA with the Khmer groups (fig. 3A) which reflects their previous connection. In addition, the high frequency of MSY haplogroup O2a\* and C\* (fig. 1), close genetic relationship to many TK- and ST-speaking groups (fig. 3B) and highest MPD value for MSY (fig. 2C) indicated later extensive gene flow, promoting the paternal difference of Nyahkur from the Mon and also other AA groups.

Previous genetic studies of G6PD mutations reported a high prevalence of the Mahidol type G6PD deficiency in the Mon, Burmese, and Karen, different from Thai, Laotian, and Khmer groups exhibiting the Vientiane-type G6PD mutation (Iwai et al. 2001; Matsuoka et al. 2005; Nuchprayoon et al. 2008). Thus, both our results and previous studies indicate a close genetic relationship among Mon, Burmese, and Karen in Myanmar, suggesting a common origin or extensive gene flow. Our previous mtDNA study also revealed genetic relations between some Mon groups (MO1 and MO5) and Burmese, with both of them close to some Indian populations, whereas other Mon groups are closer to the Karen groups (MO2, MO3, and MO4) (see details in Kutanan, Kampuansai, Brunelli, et al. [2018]). In general, genetic mixing among Mon, Karen, and Myanmar might have happened before the arrival of the Mon to Thailand, whereas mixing among the Mon and central Thai would have occurred after the arrival of the Mon. However, MSY data for the Burmese are limited, and further MSY studies of populations from Myanmar are needed to confirm this scenario.

A connection between Indian groups and the Mon is suggested by South/Central Asian MSY lineages in the Mon, for example, R\*, H\*, J\*, L\*, and Q\* (fig. 1), consistent with some mtDNA lineages, for example, W3a1b, M6a1a, M30, M40a1, M45a, and I1b (Kutanan et al. 2017; Kutanan, Kampuansai, Brunelli, et al. 2018). Thus, both mtDNA and the MSY indicated contact between the ancestors of the Mon and Indian. Archaeological evidence also suggests Indian influences, for example, the symbolism on the Dvaravati coin which indicates the importance of royalty, and includes several motifs associated with Indian precedents of the first to fourth century AD (Higham and Thosarat 2012).

### Demographic Changes

Demographic expansion of Thai/Lao populations is noticeably detected in both paternal and maternal lineages at the beginning of the Holocene, ~10 ka (fig. 6). In this period, increasing and more stable temperatures might have facilitated population expansion (Wen et al. 2016). The male  $N_e$  increase during the Holocene is primarily driven by the O2a2a\* and O2a1c\* lineages (fig. 8). The Holocene expansion might thus be related to an expansion of HM paternal lineages, as O2\* (O-M122\*) is thought to have arisen at the beginning of the Holocene near Tibet (van Driem 2017).

According to this hypothesis, the bearers of this haplogroup became the progenitors of the “Yangtzean” or HM paternal lineages, and contributed this lineage to the ancient AA who carried O1b1a1a\* or O-M95\* by sharing of knowledge about rice agriculture. However, further sequencing of MSY lineages belonging to the HM populations are needed to verify this hypothesis.

During the Neolithic period, other significant expansions are observed in almost all ethnicities and many MSY haplogroups, that is, O1b1a1a1b\*, O1b1a1a1a\*, and R\* (fig. 8). Previously, it was suggested that the demographic expansion pattern in the Neolithic in SEA shows strong expansion dynamics, different characteristics than the Paleolithic expansion, and sex-specific expansion patterns, with earlier expansions in female than in male lineages. (Wen et al. 2016). The expansion signals in our results coincide with the beginning of the SEA Neolithic ~5–4.5 ka, during which farming expanded from China to SEA (Bellwood 2018). The farming technology for food production could support a higher population density than hunting–gathering, as agriculture could produce a more steady food supply, and males could avoid hunting dangerous animals; thus, effective population size would increase (Jobling et al. 2004; Yan et al. 2014). The farmer expansion ~4 ka was probably related to ancestral AA-speaking hill tribes with predominantly O-M95\* lineages that knew rice agriculture (van Driem 2017; Lipson et al. 2018; McColl et al. 2018). However, the movement of Neolithic groups from southern China to MSEA probably involved not only AA groups but also TK groups (Bellwood 2018). In our study, a Neolithic expansion signal was observed for the MSY in all studied groups, indicating a large demographic expansion and probable admixture among the ancestors of indigenous southern Chinese groups during the Neolithic period. Haplogroup R1a was previously suggested to show a similar expansion, with paternal population growth during ~6.5–4 ka observed globally (Poznik et al. 2016; Wang et al. 2016).

In addition, we found another significant expansion during the Bronze age ~2 ka that involves TK-speaking populations, reflected by some haplogroups prevalent in the TK, for example, O1a\* (fig. 8). This TK-related expansion is consistent with the strong expansion detected in the BSP of Xishuangbanna Dai (fig. 7) and corresponds with the results of a recent ancient DNA study (McColl et al. 2018). The southward expansion of the indigenous southern Chinese TK speakers to MSEA was probably driven by the Han Chinese expansion from the Yellow River basin to southern China during the Qin dynasty, starting ~2.5 ka (Bellwood 2018). The migration and expansion of prehistoric TK groups during the Bronze Age has had a profound influence on the modern Thais and Laotians in term of languages and genes. Nowadays, TK languages are mostly concentrated in present-day Thailand and Laos, and the relatively high level of TK genetic homogeneity might be also driven by this recent expansion.

Our previous mtDNA modeling to explore the migration and expansion of prehistoric TK groups during the Bronze Age supported the spread of TK languages via demic diffusion

and admixture (Kutanan et al. 2017; Kutanan, Kampuansai, Brunelli, et al. 2018). Here, a similar modeling approach for the MSY data found weak support for cultural diffusion of TK languages. Although we built the model based on historical sources (supplementary fig. 6, Supplementary Material online), the models did not generate the observed variation (supplementary fig. 7, Supplementary Material online), indicating that the analyzed models do not correspond to the real paternal population history. A possible reason for this striking difference between maternal and paternal histories might be warfare. Historically, many areas of Thailand saw frequent warfare involving various TK groups ~200–500 ya (Penth 2000). As a result, forced migrations were imposed upon the losing side and men were taken captive more often than women because men could be used to strengthen the victors' armies. This could result in a different history for the TK male versus female population. More complex demographic models could therefore more accurately capture the paternal history of Thai/Lao populations.

It may be that the MSY sequences do not harbor enough information to distinguish among the different demographic scenarios. However, comparison of genetic differences ( $\Phi_{st}$  and corrected pairwise differences) among the groups used in the simulations does support a real contrast in the maternal versus paternal histories for the major TK groups in each region, and also finds genetic heterogeneity among these major groups. The northern Thai people showed closer genetic relationship with the Dai than AA groups in both mtDNA and MSY, supporting the demic diffusion model, whereas the ethnic Lao are closer to Dai for mtDNA but for MSY they are related to both Dai and AA rather equally, suggesting demic diffusion for the maternal history and admixture for the paternal history. The central Thai MSY sequences could be of AA origin because they are genetically more similar to the AA groups than the Dai, supporting cultural diffusion, but for mtDNA they are related to both Dai and AA rather equally, supporting admixture in central Thailand as found previously (Kutanan, Kampuansai, Brunelli, et al. 2018). Overall, these results suggest that the demographic history of Khon Mueang, ethnic Lao, and central Thais are different, possibly reflecting either different migration routes or different small TK groups that expanded from China (Higham and Thosarat 2012). In addition, different patterns of admixture for males versus females could have occurred in ethnic Lao and central Thais. Archaeological and historical evidence indicate that prior to the TK migration, there were existing rich civilizations in the area, for example, the Dvaravati of the Mon and Chenla of the old Khmer. With the arrival of TK groups, the Mon people were incorporated by intermarriage into Tai society and adopted the increasing dominant Thai language as their own (Higham and Thosarat 2012). Our results suggest that there was variation in the pattern of cultural diffusion/admixture involving males versus females in different groups in the area of northeastern and central Thailand and Laos. Such admixture could also have had an impact on the patterns of genetic diversity in the matrilineal versus patrilineal groups, which might then contribute to diminishing the genetic signal attributable to residence pattern.

Finally, another more recent expansion signal was detected in the northern Thai AA-speaking Lawa, involving haplogroups O2a2b\* and N\* (figs. 6 and 8). Historical evidence indicates that after the arrival of the TK groups in northern Thailand, the native Lawa groups were fragmented and moved to the mountains (Penth 2000), resulting in cultural and geographical isolation. In support of this model of isolation and drift, we note that the most negative Tajima's *D* value is observed in the LW3 group, which suggests population expansion after a bottleneck (fig. 2D).

## Conclusion

We compared high-resolution mtDNA and MSY sequences and found contrasts in the maternal and paternal genetic history of various Thai/Lao groups, in particular the hill tribes, the major TK groups in different regions, and the AA- and ST-speaking groups, as well as significant genetic heterogeneity among samples from the same ethnolinguistic group from different locations (figs. 1 and 4). These contrasting patterns reflect the influence of different factors in different Thai/Lao groups, for example, cultural practices in the hill tribes coupled with genetic drift in some population, as well as gene flow in the lowland Mon and TK groups. This new MSY study from Thai/Lao males provides more insight into the past demographic history in the paternal line and, along with our previous mtDNA studies, is generally in agreement with recent ancient DNA studies in SEA that indicate two demographic expansions from southern China to MSEA, with the first involving the ancestors of AA groups and the second involving TK groups (Lipson et al. 2018; McColl et al. 2018). Overall, the contrasting results for the maternal versus paternal history of some Thai/Lao groups supports the importance of detailed studies of uniparental markers, as such contrasts would not have been revealed by studying autosomal markers in just a few Thai/Lao groups. Additional ancient DNA studies, coupled with more detailed genome-wide data from present-day populations, will provide a complete reconstruction of the genetic history of this region.

## Materials and Methods

### Studied Populations

Genomic DNA was extracted from blood, buccal swab or saliva of 914 males belonging to 57 populations that were classified into 26 ethnolinguistic groups, as described previously (Kutanan et al. 2017; Kutanan, Kampuansai, Changmai, et al. 2018) (fig. 1 and supplementary table 1, Supplementary Material online). Ethical approval for this study was provided by Khon Kaen University, Narueuan University, and the Ethics Commission of the University of Leipzig Medical Faculty.

### MSY Sequences

We prepared genomic libraries for each sample using a double index scheme (Kircher et al. 2012) and enriched the libraries for ~2.34 mB of the MSY via in-solution hybridization-capture using a previously designed probe set (Kutanan, Kampuansai, Brunelli, et al. 2018) and the Agilent Sure

Select system (Agilent, CA); further details on the probe design are provided in [supplementary table 6, Supplementary Material](#) online. Sequencing was carried out on the Illumina HiSeq 2500 platform with paired-end reads of 125-bp length. Standard Illumina base-calling was performed using Bustard. Illumina adapters were trimmed and completely overlapping paired sequences were merged using leeHOM (Renaud et al. 2014). Demultiplexing of the pooled sequencing data was done by deML (Renaud et al. 2015). The alignment and post-processing pipeline of the sequencing data was described previously (Kutanan, Kampuansai, Brunelli, et al. 2018).

## Statistical Analysis

### *Genetic Diversity and Structure*

We combined the 914 newly generated sequences together with 14 published sequences (Kutanan, Kampuansai, Brunelli, et al. 2018) belong to two hunter–gatherer populations from Thailand: Mlabri and Maniq ([supplementary table 1, Supplementary Material](#) online). This study thus includes 928 MSY sequences from 59 populations and 28 ethnolinguistic groups of Thailand and Laos. To compare with the MSY data, we selected 1,434 mtDNA sequences from the same populations from our previous studies (Kutanan et al. 2017; Kutanan, Kampuansai, Brunelli, et al. 2018; Kutanan, Kampuansai, Changmai, et al. 2018) ([supplementary table 1, Supplementary Material](#) online). We used Arlequin 3.5.1.3 (Excoffier and Lischer 2010) for the following analyses: summary statistics of genetic diversity within populations, the matrix of genetic distances ( $\Phi_{st}$ ), AMOVAs, and Mantel tests of the correlation between genetic and geographic distances.

### *Genetic Relationships*

To investigate the paternal relatedness between populations, we performed a DAPC (Jombart et al. 2010). We grouped our samples based on population sampled, geographic location, and ethnicity ([supplementary table 1, Supplementary Material](#) online) before running the analysis for 100,000 iterations using *adegenet* 1.3-1 (Jombart 2008).

A correspondence analysis based on MSY haplogroup counts was performed using STATISTICA 13.0 (StatSoft, Inc., USA). Haplogroup assignment was performed by yHaplo (Poznik 2016). The R package (R Development Core Team 2016) was used to carry out a nonparametric MDS analysis (based on  $\Phi_{st}$  values of MSY and mtDNA), the MDS heat plot with five dimensions, showing per-dimension standardized values between 0 and 1, and heat plots of the  $\Phi_{st}$  distance matrix and the matrix of shared haplotypes.

To get a broad picture of population relationships in Asia, we included 552 MSY sequences from Asian groups for comparison. We downloaded the published Y chromosome sequencing data from the SGPDP data set ([https://sharehost.hms.harvard.edu/genetics/reich\\_lab/sgdp/Y-bams/Y.tar](https://sharehost.hms.harvard.edu/genetics/reich_lab/sgdp/Y-bams/Y.tar); last accessed June 25, 2018) (Mallick et al. 2016), the 1000 Genomes Project (1000 Genomes Project Consortium et al. 2015) and the study of Poznik et al. (2016). We merged and processed all sequencing data through the same pipeline as

the samples in our study (Kutanan, Kampuansai, Brunelli, et al. 2018). The resulting variant file was merged with data from previous study (Karmin et al. 2015; <http://evolbio.ut.ee/chrY/>; last accessed June 25, 2018) using Heffalump v0.2 (<https://bitbucket.org/ustenzel/heffalump>; last accessed June 25, 2018). We subset the variant file to sites that were overlapping the regions present on our capture bait and to samples that had a major haplogroup that was also present in our data set. These samples were combined with our samples; we then removed variant sites for which <25% of the samples had genotype information, and samples that had >25% of all sites with missing genotype information. The resulting data set provides 16,684 variable sites, which was imputed using BEAGLE v4.1 (Browning and Browning 2016). Additional details on these populations are provided in [supplementary table 7, Supplementary Material](#) online.

### *Bayesian Skyline Plots*

Based on Bayesian Markov Chain Monte Carlo analyses, BEAST 1.8.4 was used to construct BSPs by ethnicity and by haplogroup (Drummond et al. 2012). To avoid a false detection of bottlenecks stemming from the sample collection procedure (Heller et al. 2013), we pooled all populations within the same ethnicity and ran jModel test 2.1.7 (Darriba et al. 2012) to select the most suitable model for each run during the creation of the input file for BEAST via BEAUTi v1.8.2. We used an MSY mutation rate of  $8.71 \times 10^{-10}$  substitutions/bp/year (Helgason et al. 2015), and the BEAST input files were modified by an in-house script to add in the invariant sites found in our data set. Both strict and log normal relaxed clock models were run for each ethnicity and haplogroup, with marginal likelihood estimation (Baele et al. 2012, 2013). After each BEAST run, the Bayes factor was computed from the log marginal likelihood of both models to choose the best-fitting BSP. Tracer 1.5.0 was used to check the results. We also performed the BSP of compared populations, that is, Dai, Kinh, Southern Han, Northern Han, and Japanese from published MSY sequences (Poznik et al. 2016). The BSPs by ethnicity based on mtDNA genomes were carried out in a previous study (Kutanan, Kampuansai, Changmai, et al. 2018).

### *Approximate Bayesian Computation*

In order to investigate the paternal origin of TK groups in Thailand/Laos and their local histories, we employed five data sets (encompassing northern Thailand, central Thailand, and northeastern Thailand and Laos) and compared three competing scenarios: demic diffusion (i.e., a migration of people from southern China, who are then the ancestors of present-day Thai/Lao TK people); cultural diffusion (i.e., the Thai ancestors were the native AA groups who shifted languages and culture to TK) and continuous migration (i.e., gene flow between a migrant TK and native AA groups) that were developed based on known historical hypotheses ([supplementary fig. 6, Supplementary Material](#) online). The immigrant and endogenous scenarios postulated an initial split of AA

and Dai populations, with a subsequent treelike split of the target group from Dai (immigrant) or AA (endogenous) populations. The continuous migration model not only shared the same demographic history as the immigrant model but also allowed subsequent bidirectional migration between the newly originated population and the AA population. All of the simulations assumed uniform population sizes, fixed separation times based on historical records, a fixed mutation rate of  $8.71 \times 10^{-10}$  substitutions/bp/year (Helgason et al. 2015), and a prior distribution for both effective population sizes and migration rates (supplementary table 4, Supplementary Material online). Finally, due to the uneven sample size between the tested groups, we simulated a number of individuals equal to the lowest sample size among the populations in the model.

We simulated the derived site frequency spectrum (unfolded-SFS) for 2,364,048 loci using the fastsimcoal simulator (Excoffier and Foll 2011) with the flag -s, through the software package ABCtoolbox (Wegmann et al. 2011) and running 50,000 simulations for each model. The observed SFS was calculated with the software 4P (Benazzo et al. 2015). To determine the best performing scenario in each set we employed the model selection procedure ABC-RF (Pudlo et al. 2016), which relies on random forest machine learning methodology (Breiman 2001). This classification algorithm is trained on a reference table of simulations and allows the prediction of the most suitable model at each value of a set of covariates (i.e., the summary statistics). Additional details concerning the ABC-RF analyses are described in our previous study (Kutanan, Kampuansai, Brunelli, et al. 2018).

## Data Availability

All reads that aligned to the region of the MSY that was targeted by the capture-enrichment array were deposited in the European Nucleotide Archive (ENA) (study ID: PRJEB31636).

## Supplementary Material

Supplementary data are available at *Molecular Biology and Evolution* online.

## Acknowledgments

We would like to thank all sample donors, village chief, and coordinators, Suparat Srithawong, Sukhum Ruangchai, Khamnikone Sipaseuth, Worasitikulya Taratima, Saksuriya Triyarach, Narongdech Mahasirikul, Supada Khonyoung, Dusit Boonmekam, Tharanat Hin-on, Kantaphon Chueahor, Pittayawat Pittayaporn, and Waraporn Hongsaphinan, for assistance in collecting samples. We thank Murray Cox, Brigitte Pakendorf, and Rasmi Shoocongdej for valuable discussion. This study was supported by the Max Planck Institute for Evolutionary Anthropology. W.K. was also funded by the Thailand Research Fund (grant number RSA6180058), Khon Kaen University (grant number 6100100), KKU's Thai Visiting Scholar 2018, and the Research and Academic Affairs Promotion Fund (RAAPF), Faculty of Science, Khon Kaen University, Fiscal year 2018. M.Sr. was funded by Naresuan

University (grant number R2561B029). J.K. was funded by Chiang Mai University.

## Author Contributions

W.K. and M.S. conceived and designed the project; W.K., J.K., and M.Sr. collected samples; W.K. and R.S. generated data; W.K., A.B., S.G., L.A., A.H., and E.M. involved data analyses; W.K. and M.S. wrote the article with input from all coauthors.

## References

- 1000 Genomes Project Consortium, Auton A, Brooks LD, Durbin RM, Garrison EP, Kang HM, Korbel JO, Marchini JL, McCarthy S, McVean GA, et al. 2015. A global reference for human genetic variation. *Nature* 526(7571):68–74.
- Arias L, Schröder R, Hübner A, Barreto G, Stoneking M, Pakendorf B. 2018. Cultural innovations influence patterns of genetic diversity in Northwestern Amazonia. *Mol Biol Evol.* 35(11):2719–2735.
- Bae CJ, Douka K, Petraglia MD. 2017. Human colonization of Asia in the Late Pleistocene: an introduction to supplement 17. *Curr Anthropol.* 58(S17):S373–S382.
- Baele G, Lemey P, Bedford T, Rambaut A, Suchard MA, Alekseyenko AV. 2012. Improving the accuracy of demographic and molecular clock model comparison while accommodating phylogenetic uncertainty. *Mol Biol Evol.* 29(9):2157–2167.
- Baele G, Li WLS, Drummond AJ, Suchard MA, Lemey P. 2013. Accurate model selection of relaxed molecular clocks in Bayesian phylogenetics. *Mol Biol Evol.* 30(2):239–243.
- Bajic V, Barbieri C, Hübner A, Güldemann T, Naumann C, Gerlach L, Berthold F, Nakagawa H, Mpoloka SW, Roewer L, et al. 2018. Genetic structure and sex-biased gene flow in the history of southern African populations. *Am J Phys Anthropol.* 167(3):656–671.
- Barbieri C, Hübner A, Macholdt E, Ni S, Lippold S, Schröder R, Mpoloka SW, Purps J, Roewer L, Stoneking M, et al. 2016. Refining the Y chromosome phylogeny with southern African sequences. *Hum Genet.* 135(5):541–553.
- Bellwood P. 2018. The search for ancient DNA heads east. *Science* 361(6397):31–32.
- Benazzo A, Panziera A, Bertorelle G. 2015. 4P: fast computing of population genetics statistics from large DNA polymorphism panels. *Ecol Evol.* 5(1):172–175.
- Besaggio D, Fuselli S, Srikumool M, Kampuansai J, Castrì L, Tyler-Smith C, Seielstad M, Kangwanpong D, Bertorelle G. 2007. Genetic variation in Northern Thailand Hill Tribes: origins and relationships with social structure and linguistic differences. *BMC Evol Biol.* 7(Suppl 2):S12.
- Breiman L. 2001. Random forests. *Mach Learn.* 45(1):5–32.
- Browning BL, Browning SR. 2016. Genotype imputation with millions of reference samples. *Am J Hum Genet.* 98(1):116–126.
- Brunelli A, Kampuansai J, Seielstad M, Lomthaisong K, Kangwanpong D, Ghirotto S, Kutanan W. 2017. Y chromosomal evidence on the origin of northern Thai people. *PLoS One* 12(7):e0181935.
- Cai X, Qin Z, Wen B, Xu S, Wang Y, Lu Y, Wei L, Wang C, Li S, Huang X, et al. 2011. Human migration through bottlenecks from Southeast Asia into East Asia during last glacial maximum revealed by Y chromosomes. *PLoS One* 6(8):e24282.
- Chaix R, Quintana-Murci L, Hegay T, Hammer MF, Mobasher Z, Austerlitz F, Heyer E. 2007. From social to genetic structures in central Asia. *Curr Biol.* 17(1):43–48.
- Darriba D, Taboada GL, Doallo R, Posada D. 2012. jModelTest 2: more models, new heuristics and parallel computing. *Nat Methods.* 9(8):772.
- Demeter F, Shackelford LL, Bacon AM, Durringer P, Westaway K, Sayavongkhamdy T, Braga J, Sichanthongtip P, Khamdalavong P, Ponche JL, et al. 2012. Anatomically modern human in Southeast Asia (Laos) by 46 ka. *Proc Natl Acad Sci U S A.* 109(36):14375–14380.



- Diffloth G. 1984. The Dvaravati Old Mon Language and Nyah Kur. Bangkok (Thailand): Chulalongkorn University Print House.
- Drummond AJ, Suchard MA, Xie D, Rambaut A. 2012. A Bayesian phylogenetics with BEAUti and the BEAST 1.7. *Mol Biol Evol.* 29(8):1969–1973.
- Excoffier L, Foll M. 2011. Fastsimcoal: a continuous-time coalescent simulator of genomic diversity under arbitrarily complex evolutionary scenarios. *Bioinformatics* 27(9):1332–1334.
- Excoffier L, Lischer H. 2010. Arlequin suite ver 3.5: a new series of programs to perform population genetics analyses under Linux and Windows. *Mol Ecol Resour.* 10(3):564–567.
- Guillou E. 1999. The Mons: a civilization of Southeast Asia. Bangkok (Thailand): Siam Society Under Royal Patronage.
- Hallast P, Batini C, Zadik D, Maisano Delsler P, Wetton JH, Arroyo-Pardo E, Cavalleri GL, de Knijff P, Destro Bisol G, Dupuy BM, et al. 2015. The Y-chromosome tree bursts into leaf: 13,000 high-confidence SNPs covering the majority of known clades. *Mol Biol Evol.* 32(3):661–673.
- Helgason A, Einarsson AW, Gumundsdóttir VB, Sigursson A, Gunnarsdóttir ED, Jagadeesan A, Ebenesersdóttir SS, Kong A, Stefánsson K. 2015. The Y-chromosome point mutation rate in humans. *Nat Genet.* 47(5):453–457.
- Heller R, Chikhi L, Siegmund HR. 2013. The confounding effect of population structure on Bayesian skyline plot inferences of demographic history. *PLoS One* 8(5):e62992.
- Higham C. 2013. Hunter-gatherers in Southeast Asia: from prehistory to the present. *Hum Biol.* 85(1-3):21–43.
- Higham C. 2014. Early Mainland Southeast Asia: from first humans to Angkor. Bangkok (Thailand): River Books Press.
- Higham C. 2017. First farmers in Mainland Southeast Asia. *JIPA* 41:13–21.
- Higham C, Thosarat R. 2012. Early Thailand from prehistory to Sukhothai. Bangkok (Thailand): River Books.
- Iwai K, Hirano A, Matsuoka H, Kawamoto F, Horie T, Lin K, Tantular IS, Dachlan YP, Notopuro H, Hidayah NI, et al. 2001. Distribution of glucose-6-phosphate dehydrogenase mutations in Southeast Asia. *Hum Genet.* 108(6):445–449.
- Jobling M, Hollox E, Kivisild T, Tyler-Smith C. 2004. Agricultural expansions. In: Human evolutionary genetics. New York: Garland Publishing.
- Jombart T. 2008. adegenet: a R package for the multivariate analysis of genetic markers. *Bioinformatics* 24(11):1403–1405.
- Jombart T, Devillard S, Balloux F. 2010. Discriminant analysis of principal components: a new method for the analysis of genetically structured populations. *BMC Genet.* 11(1):94.
- Kampuansai J, Völgyi A, Kutanan W, Kangwanpong D, Pamjav H. 2017. Autosomal STR variations reveal genetic heterogeneity in the Mon-Khmer speaking group of Northern Thailand. *Forensic Sci Int Genet.* 27:92–99.
- Karmin M, Saag L, Vicente M, Wilson Sayres MA, Järve M, Talas UG, Rootsi S, Ilumäe AM, Mägi R, Mitt M, et al. 2015. A recent bottleneck of Y chromosome diversity coincides with a global change in culture. *Genome Res.* 25(4):459–466.
- Kircher M, Sawyer S, Meyer M. 2012. Double indexing overcomes inaccuracies in multiplex sequencing on the Illumina platform. *Nucleic Acids Res.* 40(1):e3.
- Kumar V, Langstieh BT, Madhavi KV, Naidu VM, Singh HP, Biswas S, Thangaraj K, Singh L, Reddy BM. 2006. Global patterns in human mitochondrial DNA and Y-chromosome variation caused by spatial instability of the local cultural processes. *PLoS Genet.* 2(4):e53.
- Kutanan W, Kampuansai J, Brunelli A, Ghirotto S, Pittayaporn P, Ruangchai S, Schröder R, Macholdt E, Srikumool M, Kangwanpong D, et al. 2018. New insights from Thailand into the maternal genetic history of Mainland Southeast Asia. *Eur J Hum Genet.* 26(6):898–911.
- Kutanan W, Kampuansai J, Changmai P, Flegontov P, Schröder R, Macholdt E, Hübner A, Kangwanpong D, Stoneking M. 2018. Contrasting maternal and paternal genetic variation of hunter-gatherer groups in Thailand. *Sci Rep.* 8:1536.
- Kutanan W, Kampuansai J, Fuselli S, Nakbunlung S, Seielstad M, Bertorelle G, Kangwanpong D. 2011. Genetic structure of the Mon-Khmer speaking groups and their affinity to the neighbouring Tai populations in Northern Thailand. *BMC Genet.* 12:56.
- Kutanan W, Kampuansai J, Srikumool M, Kangwanpong D, Ghirotto S, Brunelli A, Stoneking M. 2017. Complete mitochondrial genomes of Thai and Lao populations indicate an ancient origin of Austroasiatic groups and demic diffusion in the spread of Tai-Kadai languages. *Hum Genet.* 136(1):85–98.
- Lipson M, Cheronet O, Mallick S, Rohland N, Oxenham M, Pietruszewski M, Pryce TO, Willis A, Matsumura H, Buckley H, et al. 2018. Ancient genomes document multiple waves of migration in Southeast Asian prehistory. *Science* 361(6397):92–95.
- Mallick S, Li H, Lipson M, Mathieson I, Gymrek M, Racimo F, Zhao M, Chennagiri N, Nordenfelt S, Tandon A, et al. 2016. The Simons genome diversity project: 300 genomes from 142 diverse populations. *Nature* 538(7624):201–206.
- Matsuoka H, Nguon C, Kanbe T, Jalloh A, Sato H, Yoshida S, Hirai M, Arai M, Socheat D, Kawamoto F. 2005. Glucose-6-phosphate dehydrogenase (G6PD) mutations in Cambodia: G6PD Viangchan (871G>A) is the most common variant in the Cambodian population. *J Hum Genet.* 50(9):468–472.
- McColl H, Racimo F, Vinner L, Demeter F, Gakuhari T, Moreno-Mayar JV, van Driem G, Gram Wilken U, Seguin-Orlando A, de la Fuente Castro C, et al. 2018. The prehistoric peopling of Southeast Asia. *Science* 361(6397):88–92.
- Nahhas RW. 2007. Sociolinguistic survey of Lawa in Thailand. Chiang Mai (Thailand): Survey Unit Department of Linguistics Faculty of Humanities Payap University.
- Nuchprayoon I, Louicharoen C, Charoenvej W. 2008. Glucose-6-phosphate dehydrogenase mutations in Mon and Burmese of southern Myanmar. *J Hum Genet.* 53(1):48–54.
- Ocharoen S. 1998. Mons in Thailand [in Thai]. Bangkok (Thailand): Thailand Research Fund.
- Oota H, Settheetham-Ishida W, Tiwawech D, Ishida T, Stoneking M. 2001. Human mtDNA and Y-chromosome variation is correlated with matrilineal versus patrilineal residence. *Nat Genet.* 29(1):20–21.
- Penth H. 2000. A brief history of Lanna: civilizations of North Thailand. Chiang Mai (Thailand): Silkworm Books.
- Penth H, Forbes A. 2004. The people of mountaintops. In: Penth H, Forbes A, editors. A brief history of Lan Na and the peoples of Chiang Mai. Chiang Mai (Thailand): Chiang Mai City Arts and Cultural Centre Chiang Mai Municipality. p. 247–254.
- Pon Nya M. 2001. Ethnic identity and political autonomy of the Mon. In: McCormick P, Jenny M, Baker C, editors. The Mon over two millennia: monuments, manuscripts, movements. Bangkok (Thailand): Institute of Asian Studies, Chulalongkorn University. p. 169–202.
- Poznik GD. 2016. Identifying Y-chromosome haplogroups in arbitrarily large samples of sequenced or genotyped men, unpublished data. Available from: <https://www.biorxiv.org/content/early/2016/11/19/088716> (last accessed May 8, 2018).
- Poznik GD, Xue Y, Mendez FL, Willems TF, Massaia A, Wilson Sayres MA, Ayub Q, McCarthy SA, Narechania A, Kashin S, et al. 2016. Punctuated bursts in human male demography inferred from 1,244 worldwide Y-chromosome sequences. *Nat Genet.* 48(6):593–599.
- Pudlo P, Marin JM, Estoup A, Cornuet JM, Gautier M, Robert CP. 2016. Reliable ABC model choice via random forests. *Bioinformatics* 32(6):859–866.
- R Development Core Team. 2016. R: a language and environment for statistical computing. Vienna (Austria): R Foundation for Statistical Computing. Available from: <http://www.R-project.org/>; last accessed May 8, 2017.
- Renaud G, Stenzel U, Kelso J. 2014. LeeHom: adaptor trimming and merging for Illumina sequencing reads. *Nucleic Acids Res.* 42(18):e141.
- Renaud G, Stenzel U, Maricic T, Wiebe V, Kelso J. 2015. deML: robust demultiplexing of Illumina sequences using a likelihood-based approach. *Bioinformatics* 31(5):770–772.
- Saraya D. 1999. (Sri) Dvaravati: the initial phase of Siam's history. Bangkok (Thailand): Muang Boran Publishing House.

- Schliesinger J. 2000. Ethnic groups of Thailand: non-Tai-speaking peoples. Bangkok (Thailand): White Lotus Press.
- Schliesinger J. 2001. Tai group of Thailand. Bangkok (Thailand): White Lotus Press.
- Shoocongdej R. 2006. Late Pleistocene activities at the Tham Lod rock-shelter in highland Pang Mapha, Mae Hongson Province, Northwestern Thailand. In: Bacus EA, Glover IC, Pigott VC, editors. *Uncovering Southeast Asia's past*. Singapore: NUS Press. p. 22–37.
- Simons GF, Fennig CD. 2018. *Ethnologue: languages of the World*. 21st ed. Dallas (TX): SIL International.
- Sun H, Zhou C, Huang X, Lin K, Shi L, Yu L, Liu S, Chu J, Yang Z. 2013. Autosomal STRs provide genetic evidence for the hypothesis that Tai people originate from Southern China. *PLoS One* 8(4):e60822.
- van Driem GL. 2017. The domestications and the domesticators of Asian rice. In: Robbeets M, Savelyev A, editors. *Language dispersal beyond farming*. Amsterdam: John Benjamins Publishing Company. p. 183–214.
- Wang CC, Huang Y, Yu X, Chen C, Jin L, Li H. 2016. Agriculture driving male expansion in Neolithic Time. *Sci China Life Sci* 59(6):643–646.
- Wegmann D, Leuenberger C, Neuenschwander S, Excoffier L. 2011. ABCtoolbox: a versatile toolkit for approximate Bayesian computations. *BMC Bioinformatics* 11(1):116.
- Wen S-Q, Tong X-Z, Li H. 2016. Y-chromosome-based genetic pattern in East Asia affected by Neolithic transition. *Quat Int*. 426:50–55.
- Wilkins JF, Marlowe FW. 2006. Sex-biased migration in humans: what should we expect from genetic data? *Bioessays* 28(3):290–300.
- Xu S, Kangwanpong D, Seielstad M, Srikummool M, Kampuansai J, Jin L, Consortium THP-AS. 2010. HUGO Pan-Asian SNP Consortium. Genetic evidence supports linguistic affinity of Mlabri-a hunter-gatherer group in Thailand. *BMC Genet*. 11(1):18.
- Yan S, Wang C-C, Zheng H-X, Wang W, Qin Z-D, Wei L-H, Wang Y, Pan X-D, Fu W-Q, He Y-G, et al. 2014. Y chromosomes of 40% Chinese descend from three Neolithic super-grandfathers. *PLoS One* 9(8):e105691.