

SCIENTIFIC REPORTS



OPEN

Genetic and chemical differentiation characterizes top-geoherb and non-top-geoherb areas in the TCM herb rhubarb

Xumei Wang¹, Li Feng¹, Tao Zhou¹, Markus Ruhsam², Lei Huang³, Xiaoqi Hou³, Xiaojie Sun¹, Kai Fan¹, Min Huang¹, Yun Zhou¹ & Jie Song¹

Medicinal herbs of high quality and with significant clinical effects have been designated as top-geoherbs in traditional Chinese medicine (TCM). However, the validity of this concept using genetic markers has not been widely tested. In this study, we investigated the genetic variation within the *Rheum palmatum* complex (rhubarb), an important herbal remedy in TCM, using a phylogeographic (six chloroplast DNA regions, five nuclear DNA regions, and 14 nuclear microsatellite loci) and a chemical approach (anthraquinone content). Genetic and chemical data identified two distinct groups in the 38 analysed populations from the *R. palmatum* complex which geographically coincide with the traditional top-geoherb and non-top-geoherb areas of rhubarb. Molecular dating suggests that the two groups diverged in the Quaternary c. 2.0 million years ago, a time of repeated climate changes and uplift of the Qinghai-Tibetan Plateau. Our results show that the ancient TCM concept of top-geoherb and non-top-geoherb areas corresponds to genetically and chemically differentiated groups in rhubarb.

Traditional Chinese Medicine (TCM) has developed over millenia in China and has exerted its influence on medical culture in Asia for more than a thousand years. Many unique concepts have formed throughout this long historical process, one of which is the concept of geo-herbalism. The idea behind geo-herbalism is that the medicinal properties of herbs depend on the specific area of collection (e.g., growing on certain mountains, valleys, or in a specific province) and that plants collected from some areas have higher medicinal value than those collected from other areas. The ancient practitioners of Chinese medicine called medicinal plants obtained from regions with allegedly superior efficacy top-geo-herbs, while the ones from other regions were known as non-top-geoherbs^{1–3}, suggesting a profound understanding of the differentiation between or within plant species.

The mentioning of top-geoherbs and the importance of provenance for medicinal plants in Chinese medicine was first discussed in *Shen Nong's Herbal*, which was written during the Eastern Han Dynasty (25–220 C.E.). Among the 500 commonly used TCMs, there are about 200 TCMs where the concept of geo-herbalism is applicable⁴. The consumption of these 200 TCMs accounts for 80% of the total consumption of TCMs^{5,6}.

It is well known that species are the fundamental units in all fields of biology^{7–10} and species differentiation can be demonstrated by morphological, chemical, karyological, and genetic means. The first three traits are phenotypic in nature and might not necessarily reveal species differentiation because of phenotypic plasticity¹¹. It is believed that plant secondary metabolites (PSMs) are the important constituents in medicinal plants, and differences in PSMs between top-geoherbs and non-top-geoherbs might not be the presence or absence of one or more compounds, but rather the content variation or the unique combination of certain compounds³. Over the last century, the differences in the composition of the main medicinal compounds between top- and non-top-geoherbs have been tested through various phytochemistry studies^{3,5,12,13}. For example, the content of sarsasapogenin in *zhi mu* (*Rhizoma Anemarrhenae*) from top-geoherb regions was three times higher than that in non-top-geoherb regions¹³; the content of phenols in top-geoherb regions of *hou po* (*Cortex Magnoliae Officinalis*) was six times higher than that in non-top-geoherb regions; also see *jin yin hua* (*Flos Lonicerae*) and *chuan xin lian* (*Herba*

¹School of Pharmacy, Xi'an Jiaotong University, Xi'an, 710061, China. ²Royal Botanic Garden Edinburgh, Edinburgh, EH3 5LR, UK. ³College of Life Sciences, Shaanxi Normal University, Xi'an, 710062, China. Xumei Wang, Li Feng and Tao Zhou contributed equally to this work. Correspondence and requests for materials should be addressed to X.W. (email: wangxumei@mail.xjtu.edu.cn) or M.R. (email: MRuhsam@rbge.org.uk)

Andrographitis Paniculatae^{3,5,12}. Although few studies have attempted to explain the genetic variation between top-geoherb and non-top-geoherb (but see Guo *et al.*¹⁴ and Yuan *et al.*¹⁵), the abovementioned studies were conducted either on very limited populations or medicinal plants derived from cultivars. In addition, whether chemical differentiation also corresponds to genetic differentiation between top-geoherb and non-top-geoherb has not yet been investigated. Phylogeography, the analysis of the spatio-temporal pattern of population genetic variation, has been used to elucidate the evolutionary history and differentiation of populations within species^{16,17} which is often initiated by large scale events such as the climatic fluctuations observed during the Quaternary^{16–20}.

Rhubarb, known as the “lord or king of herbs” in China, is an important TCM that has been used for over 2,000 years as a purgative medicine. The herbal remedy “rhubarb” consists of the dried roots and rhizomes of any species of *Rheum officinale* Baill., *R. palmatum* Linn., or *R. tanguticum* (Maxim. ex Regel) Maxim. ex Balf. (Polygonaceae). The major medicinally active ingredients in rhubarb are anthraquinones, which were used to assess the quality of rhubarb in the *Pharmacopoeia of the People's Republic of China 2015*. All three species are endemic to China and form a clade based on molecular data²¹. According to *Flora of China*, the delimitation of these three rhubarb species is based primarily on the degree of leaf blade dissection and the shape of the lobes²². The blades of *R. officinale* and *R. palmatum* are lobed; the lobes of *R. officinale* are broadly triangular, and those of *R. palmatum* are narrowly triangular. The blades of *R. tanguticum* are parted, and its lobes are narrow and triangular-lanceolate²². However, the results of a morphological analysis indicated that the degree of leaf blade dissection is continuous among the three species¹⁷. Geographically, *R. palmatum* is the most widely distributed species from the eastern margin of the Qinghai-Tibetan Plateau (QTP) to Qinling extending to the Zhongtiao Mountains (Shanxi Province). The northwestern and southeastern distributions of *R. tanguticum* and *R. officinale* partly overlap with *R. palmatum*. The three species co-occur in northwest Sichuan Province, and form mixed populations²³. Molecular evidence using inter-simple sequence repeats (ISSR) markers have shown that the populations from each species did not cluster according to species^{23,24}. Genetic and morphological data therefore indicate that these three species can be regarded as one species (henceforth named the *R. palmatum* complex). In TCM, rhubarb plants collected from Qinghai, Gansu, and Sichuan Provinces (in or near the QTP and the Hengduan Mountains) are classed as top-geoherb compared to rhubarb collected further east, which are regarded to be non-topgeoherb with inferior quality. Therefore, it is likely that chemical and possibly genetic differentiation exists in the *R. palmatum* complex.

In the present study, we investigate whether chemical and genetic differences can be detected between top-geoherb and non-top-geoherb areas of the *R. palmatum* complex throughout its distribution in China. The specific objectives are to: (1) quantify the anthraquinone content from top-geoherb and non-top-geoherb areas; (2) determine whether genetic differentiation exists between top-geoherb and non-top-geoherb areas of rhubarb; (3) investigate the possible reasons that might have resulted in the genetic variation between top-geoherb and non-top-geoherb areas of rhubarb; and (4) test whether genetic differentiation is consistent with the concept of top-geoherb and non-top-geoherb in TCM.

Results

Analysis of anthraquinone content. The HPLC results showed that eight anthraquinones were successfully detected in all samples except for aloe-emodin-8-O- β -D-glucopyranoside and rhein-8-O- β -D-glucopyranoside, which could not be found in any sample (Fig. 1a; Supplementary Fig. S1 and Table S1). Emodin presented the greatest difference, with the highest content 55.00 times the lowest content. Chrysophanol-8-O- β -D-glucopyranoside showed the least difference, with the highest content 2.94 times the lowest content. PCA results indicated that the samples could be divided into two clusters, an eastern and a western group (Supplementary Fig. S2). The content of the three constituents (rhein, emodin, and emodin-8-O- β -D-glucopyranoside) in the western group was significantly higher than in the eastern group (t-test, $P < 0.05$), while the physcion-8-O- β -D-glucopyranoside in the eastern group was significantly higher than that in western group ($P < 0.05$) (Fig. 1b). We also determined similar results when compared the differences of chemical composition between the western and eastern clades using general liner model.

Sequence variation and population structure. Sequencing of the six cpDNA segments from 377 individuals from the *R. palmatum* complex resulted in a joint alignment length of 3,316 bp and 34 haplotypes defined by 57 substitutions and 14 indels. Populations were fixed for a single haplotype (Supplementary Fig. S3 and Table S2), except for five populations (GSZQ, HBXS, SCKD, SNH, and SNZZ), and the cpDNA gene tree didn't show meaningful structure that the nuclear genes and nSSR data did (see below). The total haplotype (Hd) and nucleotide diversities (π) were 0.960 and 0.0031, respectively. Non-hierarchical analysis of molecular variance (AMOVA) indicated that the genetic variation was mainly distributed among populations rather than within populations (98.69% vs. 1.31%, Table 1). Significant phylogeographic structures was detected ($N_{ST} = 0.987$, $G_{ST} = 0.938$, $P < 0.05$). Additionally, neutrality test statistics were all non-significant in combined sequences, which was also supported by the mismatch distribution analysis (multimodal; Supplementary Fig. S4) via a significant SSD value (0.020, $P = 0.043$).

The total alignment length of four nuclear genes was 3,135 bp displaying 67 variable sites and one indel ranging in length from 708 to 818 bp. The aligned sequences of *CHS* were 1,549 bp containing 68 variable sites and four indels yielding 95 haplotypes (Fig. 2d). Total haplotype diversity (Hd) among five nuclear loci ranged from $Hd = 0.477$ to 0.944, Watterson's theta (θ_{wt}) from $\theta_{wt} = 0.0026$ to 0.0071, and the silent nucleotide diversity (π_{sil}) as well as the total nucleotide diversity (π_t) ranged from $\pi_{sil} = 0.00058$ to 0.00245 and $\pi_t = 0.00143$ to 0.00505, respectively. All neutrality test statistics were nonsignificant (Supplementary Table S3). According to the STRUCTURE analysis, the log-likelihood values of the four nuclear genes and the *CHS* gene increased with K , however delta K indicated that the optimal value for K was two (Fig. 3b), clustering the samples into an eastern and western clade (Fig. 3c; Supplementary Fig. S5b). Non-hierarchical AMOVA showed strong population

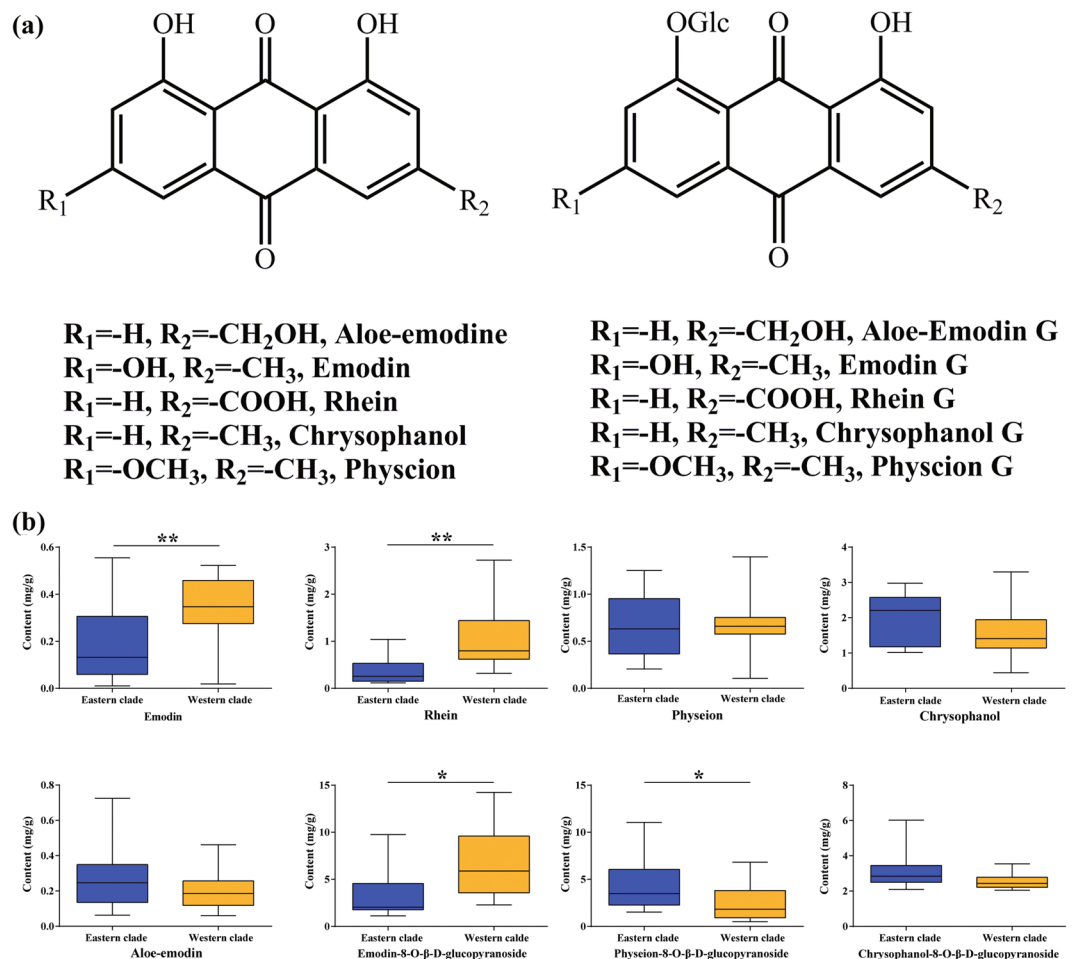


Figure 1. The chemical structures of anthraquinones and a comparative analysis. (a) The free anthraquinones and anthraquinone glycosides in rhubarb; (b) Histograms comparing the contents of chemical compounds between rhubarb samples in two clades. ** $P < 0.01$, * $P < 0.05$ (two-tailed, unpaired T-test).

structure in the *R. palmatum* complex ($\Phi_{ST} = 0.75$, $P < 0.01$). However, based on the STRUCTURE analysis, the hierarchical AMOVA indicated that 59.37% (four genes) and 57.35% (*CHS*) of the total variation was distributed among the two groups, and 22.93% (four genes) and 25.09% (*CHS*) of variation existed among populations within groups (Table 1).

Nuclear microsatellite diversity and population structure. No nuclear microsatellite locus consistently deviated from Hardy–Weinberg Equilibrium (HWE) in all the studied populations after Bonferroni corrections (corrected $\alpha = 0.00013$, $P < 0.01$). We did not find any evidence for linkage disequilibrium (LD) in any pair of nSSR loci in any population (exact tests; all $P > 0.05$), nor was there evidence for the existence of null alleles in any locus. Diversity estimates based on 14 microsatellite loci varied between all populations (Supplementary Tables S3 and S4). Population differentiation was significant with an overall $F_{ST} = 0.214$ ($P < 0.05$), and the standardized genetic differentiation G'_{ST} was higher than F_{ST} across all loci ($G'_{ST} = 0.436$, Supplementary Tables S4). AMOVA demonstrated significant genetic differentiation at the range-wide scale ($\Phi_{ST} = 0.2023$, $P < 0.01$), with 1.46% of the variation partitioned between the western and eastern clades, and 18.77% among populations within clades (Table 1). There was significant isolation-by-distance (IBD) among all populations ($r = 0.5484$, $P < 0.001$), and this effect also persisted when each sub-cluster was analyzed separately, except for sub-cluster pop4 (Supplementary Fig. S6).

The estimated values of $\text{Ln}P(K)$ increased progressively from $K = 1$ up to $K = 38$ in the STRUCTURE analyses. The highest ΔK was observed at $K = 2$ but was also significant at $K = 5$ (Fig. 4a). At $K = 2$, the cluster membership coefficient estimates suggested that clinal variation occurred along an east-west gradient separating two clusters and a geographical discontinuity arose at *ca.* 102°E. Cluster I in western China mainly distributed in the QTP and the Hengduan Mountains, and cluster II located near the Sichuan Basin and extended to the eastern areas (Fig. 4b,d). With $K = 5$, the clusters could be further subdivided into two sub-clusters (pop1, pop2) and three sub-clusters (pop3, pop4, pop5) for cluster I and II, respectively (see Fig. 4c,e). However, pop3, except for population from Miaoxian of Sichuan (SCM), was included in the western cluster based on nDNA loci. In addition, indicators of a recent bottleneck were detected in eight and five populations under the stepwise mutation model

Source of variation	d.f.	Sum of squares	Variation components	Percentage of variation	Φ -statistics
cpDNA					
All samples					
Among populations	37	3083.11	8.3216	98.69	
Within populations	342	37.900	0.1108	1.31	$\Phi_{ST} = 0.9870^{***}$
nDNA					
All samples					
Among populations	37 (34)	2671.78 (1696.28)	3.3857 (5.3098)	74.92 (75.25)	
Within populations	762 (283)	863.57 (494.30)	1.1333 (1.7466)	25.08 (24.75)	$\Phi_{ST} = 0.7492^{***}$ (0.7525**)
Two clades					
Among clades	1	1524.69 (894.07)	3.8023 (5.7030)	59.37 (57.35)	$\Phi_{CT} = 0.5937^{***}$ (0.5937***)
Among populations within clades	36 (33)	1147.08 (802.21)	1.4688 (2.4951)	22.93 (25.09)	$\Phi_{SC} = 0.5645^{***}$ (0.5645***)
Within populations	762 (283)	863.57 (494.30)	1.1333 (1.7466)	17.70 (17.56)	$\Phi_{ST} = 0.5937^{***}$ (0.8244**)
nSSR					
All samples					
Among populations	37	1455.73	0.9411	19.63	
Within populations	1396	5379.58	3.8536	80.37	$\Phi_{ST} = 0.1963^{**}$
Two clades					
Among clades	1	86.27	0.0705	1.46	$\Phi_{CT} = 0.0146^{***}$
Among populations within clades	36	1369.46	0.9069	18.77	$\Phi_{SC} = 0.1905^{***}$
Within populations	1396	5379.58	3.8536	79.77	$\Phi_{ST} = 0.2023^{***}$

Table 1. Hierarchical analysis of molecular variance (AMOVA) of samples from the *Rheum palmatum* complex. Note: Φ_{SD} differences among populations; Φ_{CT} difference among clades; Φ_{SC} , difference among populations within clades; $^{***}P < 0.001$, $^{**}P < 0.01$, $^{*}P < 0.05$, 10000 permutations; the numbers in the parentheses mean the results from *CHS* gene sequences. d.f. = Degrees of freedom.

(SMM) and the two-phase model (TPM) models, respectively. Similarly, in the mode-shift test, 11 populations suffered bottlenecks (Supplementary Table S5).

Historical and contemporary gene flow. Based on the nSSR data, all 20 pairwise estimates of historical gene flow within and between western and eastern clades were significant (no 95% confidence intervals overlapping zero), ranging from 0.23 (pop1 to pop4) to 2.88 (pop3 to pop5). The gene flow from the western clade to the eastern clade was 7.76 (95% CI: 7.25–8.32), and that of the opposite direction was 5.00 (95% CI: 4.65–5.39). Some migration rates were distinctly asymmetrical, for instance, being much greater from pop4 to pop1 (0.52) than the reverse (0.23) (Supplementary Figs S7a and b).

The mean contemporary gene flow (*mc*) from the west to east was 0.0028 (95% CI: 0–0.0071) and the opposite was 0.0021 (95% CI: 0–0.0050). The range of contemporary gene flow was 0.0016 (pop5 to pop4) to 0.0098 (pop4 to pop1). Some pairs also had asymmetrical values, including the pairs of pop2 and pop4, pop1 and pop4, and pairs within the eastern clade (Supplementary Fig. S7b).

ABC- and Bayesian skyline plot-based inferences of population history. In step 1 of the DIYABC analysis, the posterior probability for scenario 1 was 0.7213 (95% CI: 0.7124–0.7301), much higher than for the other scenarios. This scenario depicted a situation that the three sub-clusters in the eastern clade diverged simultaneously. Then in step 2, scenario 1 also had the highest posterior probability (0.5708, 95% CI: 0.5233–0.6183). This scenario indicated that the divergences of the sub-clusters in the western and eastern clades were also simultaneous. In the final step, the simulations indicated the posterior probabilities for scenarios 1 and 2 were 0.4111 (95% CI: 0.4008–0.4212) and 0.4985 (95% CI: 0.4884–0.5086), respectively, much higher than for the other two scenarios (Fig. 5).

In scenario 1 of the final step, the median values of the effective population sizes of pop1 to pop5 and NA (ancestral population) were 3.88×10^5 , 6.42×10^5 , 6.50×10^5 , 6.09×10^5 , 6.61×10^5 , and 4.66×10^5 , respectively (Table 2). The divergence times within western and eastern clades (t1), between the two clades (t2), and the time of ancient population size changes (t3) were 2.41×10^5 , 4.14×10^5 , and 1.47×10^6 generations ago, respectively. Assuming the generation time to be five years, the times of t1, t2, and t3 corresponded to 1.21, 2.07, and 6.50 million years ago (Ma), respectively. The estimated median mutation rates and the proportion of multiple step mutations were 2.84×10^{-5} and 0.568, respectively (Table 2).

In scenario 2, the effective population sizes of pop1 to pop5 and NA were 5.73×10^5 , 4.56×10^5 , 6.56×10^5 , 6.23×10^5 , 6.09×10^5 , and 3.33×10^5 , respectively. The divergence times between sub-clusters within two clades (t1), between two clades (t2), and the time of changes of the ancient population size were 2.13×10^5 , 3.68×10^5 , and 1.31×10^6 generations ago, corresponding to 1.07, 1.84, and 6.05 Ma, respectively. The estimated median mean mutation rate and the proportion of multiple step mutations were 2.89×10^{-5} and 0.471, respectively (Table 2).

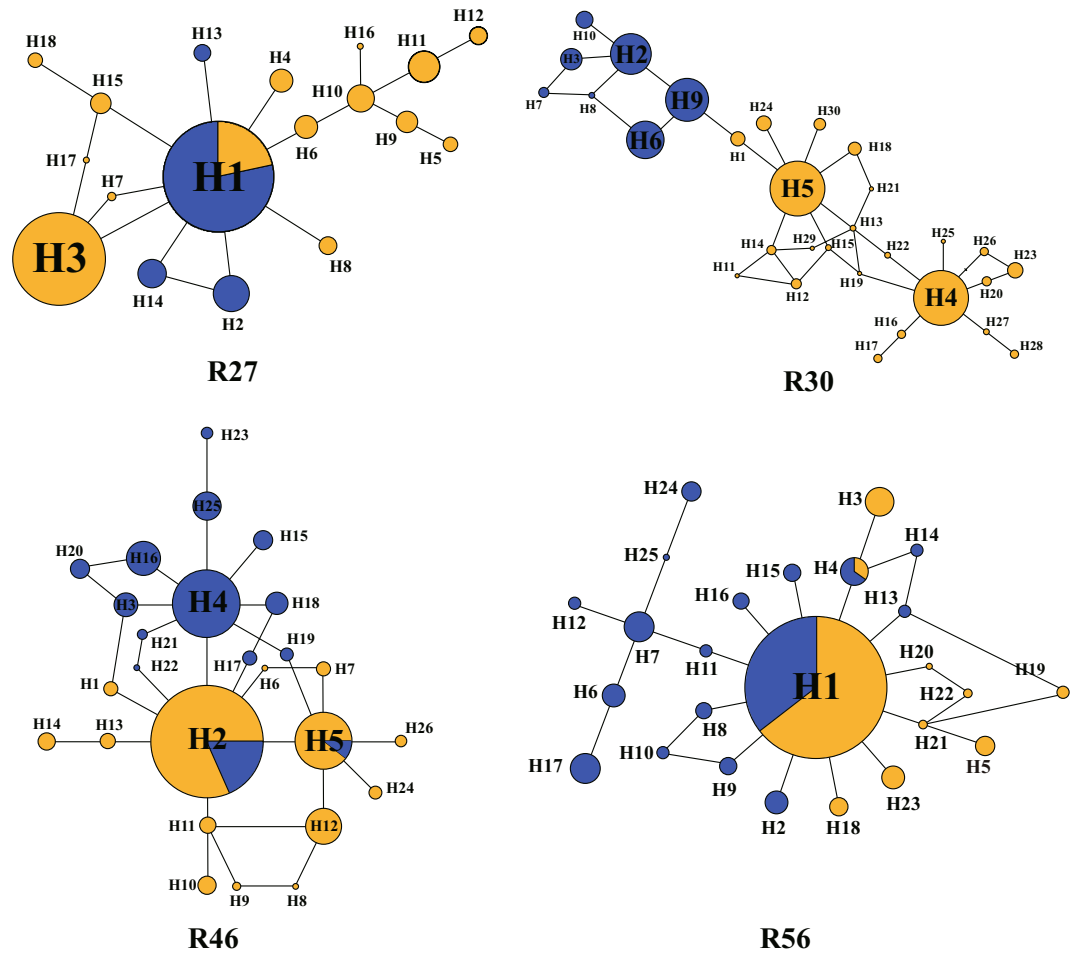


Figure 2. Haplotype genealogies for four nuclear loci. The color of the circle indicates the clade of the *Rheum palmatum* complex based on the STRUCTURE analysis of nuclear genes. The size of the circle is proportional to the frequency of the haplotype.

According to the Bayesian skyline plot (BSP) analyses, the *R. palmatum* complex experienced an expansion beginning at *ca.* 0.4 Ma and have increased its effective population size dramatically since *ca.* 0.1 Ma based on cpDNA data, whereas for the four nuclear genes each (i.e., R27, R30, R46, R56), the corresponding expansion times for the R30 gene and other three genes were all at *ca.* 0.15 Ma, despite with different expansion degrees (Supplementary Fig. S8).

Ecological niche modeling and climatic data analysis. The ecological niche models (ENMs) had a high predictive power for the two gene pools, with AUC = 0.92 and 0.93 for western and eastern clades, respectively. This is fairly congruent with their current distribution, except for some eastern clade populations that are missing in current conditions compared with the genetic architecture of this species complex (Figs 3, 4 and Supplementary Fig. S5). However, the predicted models of two clades depicted different demographic histories. For the western clade, the ENM results demonstrated a continuous expansion since the last interglacial (LIG), whereas the eastern one suggested distribution expansion from LIG to the Last Glacial Maximum (LGM) and then shrinkage from LGM to the present day (Fig. 6). Nevertheless, despite a substantial lack of precision for LGM models, a common trend was that models produced using localities from either western or eastern gene pools alone showed little overlap of their predicted distributions for all periods considered (Fig. 6a), suggesting that eastern and western clusters have occupied environmentally different regions for a substantial amount of time. Additionally, although the two models of LGM suggested that the two populations both underwent range expansions from LIG to LGM, recent studies suggested that the Interdisciplinary Research on Climate (MIROC) simulation has been shown to be more realistic than the Community Climate System Model (CCSM) in predicting potential habitats of LGM in Asia²⁵. Assuming this finding also applies to continental East Asia, we relied more on the MIROC simulation for the discussion of the projected LGM distribution. Overall, the whole distribution area of the *R. palmatum* complex had a northwestward shift from the past to the present.

Our analysis of environmental factors indicated that all 20 climatic variables contributed to the divergence between the western and eastern clades (Kruskal–Wallis tests: $P < 0.05$). The frequency distributions of each clade for each environmental variable are presented as kernel density plots (Supplementary Fig. S9), and the results from the MANOVA also distinguished significant differences between clades with regard to the environment

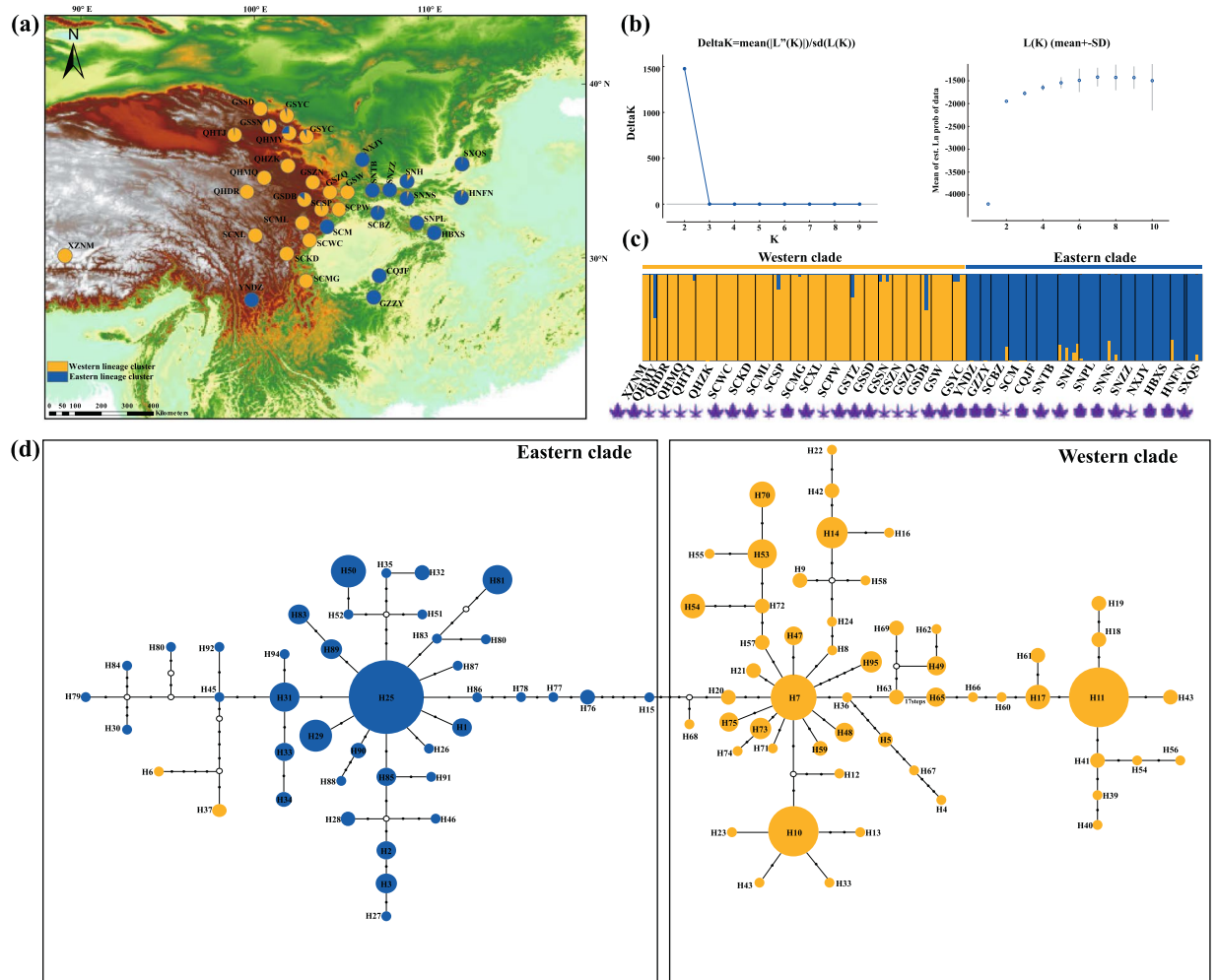


Figure 3. Bayesian clustering and haplotype network for the *CHS* gene. **(a)** Geographic origin of the 35 *Rheum palmatum* complex populations and their color-coded grouping according to the STRUCLUSTURE analysis. Population codes are identified in Table S1; **(b)** The left diagram indicates the corresponding delta K statistics calculated according to Evanno *et al.*⁶⁴, the number of clusters (K) varied from one to 10 in 10 independent runs, whereas the right plot represents changes of the mean posterior probability ($\log_e P(D)$) values of each K calculated according to Pritchard *et al.*⁶³; **(c)** Histogram of the STRUCLUSTURE analysis for the model with $K=2$ (showing the most optimal delta K). The smallest vertical bar represents one individual; **(d)** TCS-derived network of genealogical relationships between haplotypes. Each circle means a single haplotype sized in proportion to its frequency. Small open circles represent missing haplotypes. The base map was drawn using ArcGIS v.10.2 (ESRI, Redlands, CA, USA). Sketch on the bottom of the histogram of the STRUCLUSTURE analysis represents different leaf shape of each population.

occupied ($P < 0.001$). The first two PCA axes explained 82% of the variation for the present climate (Fig. 6c). However, unlike other studies on alpine plants which had loadings of different variables in PC1 and PC2^{26,27}, our PCA indicated almost all variables had high loadings for PC1 (Supplementary Table S6). Moreover, identity tests based on six and all 19 climatic variables indicated the two clades occupy identical climatic niches ($P < 0.05$) (Fig. 6b). Additionally, species records defined clear groups along the first axis of the multivariate space (Fig. 6c), indicating that substantial differences exist in the climate for each geographic region.

Discussion

We investigated whether the concept of geo-herbalism in TCM is consistent with chemical and genetic differences between top-geoh herbs and non-top-geoh herbs using the *R. palmatum* complex (rhubarb), an important herbal remedy, as a case study.

As previous findings^{23,24,28,29} showed that the three species of the *R. palmatum* complex (*Rheum officinale*, *R. palmatum*, *R. tanguticum*) do not show clear morphological separation, we used a number of molecular markers (five nDNA regions, 14 nuclear microsatellite markers and six cpDNA regions) to ascertain whether the three species can be told apart genetically. As none of the used markers retrieved three clearly separated groups (Figs 3c and 4b; Supplementary Fig. S5b) and morphological identification of the three species is unreliable, we concluded that the *R. palmatum* complex can be treated as one species.

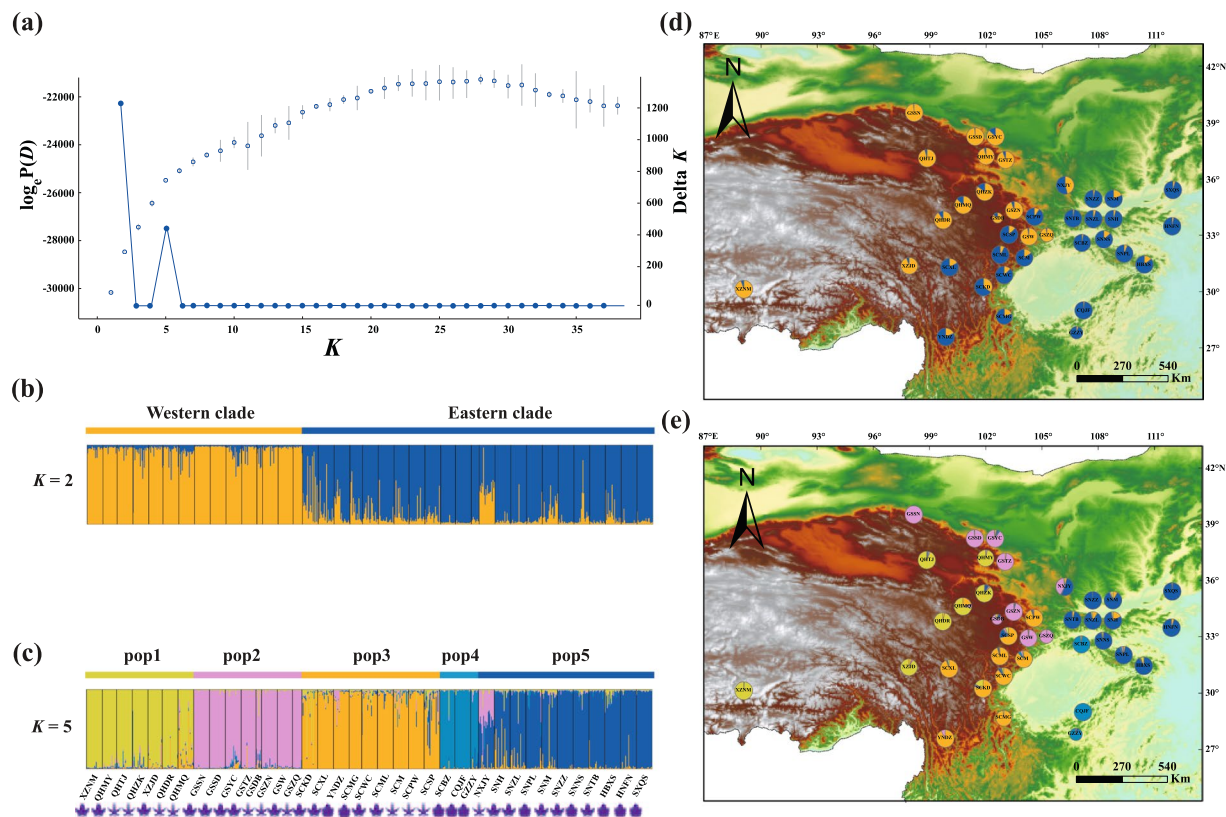


Figure 4. Bayesian clustering results of the STRUCTURE analysis for 38 populations of the *Rheum palmatum* complex for 14 SSR loci. (a) Represents changes of the mean posterior probability ($\log_e P(D)$) values of each K calculated according to Pritchard *et al.*⁶³ and the corresponding delta K statistics calculated according to Evanno *et al.*⁶⁴, respectively; (b,c) Histograms of the STRUCTURE analysis for the model with $K=2$ (b), showing the highest delta K and with $K=5$ (c, the second most optimal delta K). The smallest vertical bar represents one individual; (d,e) Geographic origin of the 38 *R. palmatum* complex populations and their color-coded grouping according to the STRUCTURE analysis when $K=2$ and $K=5$, respectively. Population codes are identified in Table S1. The base map was drawn using ArcGIS v.10.2 (ESRI, Redlands, CA, USA). Sketch on the bottom of the histogram of the STRUCTURE analysis represents different leaf shape of each population.

Analysis of the anthraquinone content, the main medicinal compound in rhubarb, of 38 populations collected from throughout the distribution range showed that the *R. palmatum* complex was divided into two geographic clusters, an eastern and a western group, which coincided with the top-geoherb and non-top-geoherb areas for rhubarb in China (Supplementary Fig. S2).

The two geographic groups were retrieved using 14 nuclear microsatellite markers suggested that the *R. palmatum* complex consists of two distinct clusters with a major phylogeographic break near the Sichuan Basin (Table 1; Figs 3c and 4b; Supplementary Fig. S5b). However, there was some disagreement between the chemical and nuclear data set regarding the placement of three populations (SNM, SCBZ, and GSYC). SNM and SCBZ genetically belonging to the eastern cluster grouped with the western clade chemically whereas GSYC grouped with the western cluster genetically but with the eastern cluster chemically. These three populations might be an exception from the close correlation between genetic relatedness and anthraquinone content because of a possible unusual combination of environmental factors which affect the chemical composition in these specific regions.

The consistent grouping of the *R. palmatum* complex into two geographic clusters suggests the survival in at least two refugia since the beginning of the Quaternary. Similar results have been reported for other temperate plants with similar ranges (e.g., *Tetracentron sinense* and *Ligularia hodgsonii*)^{30,31}. These divergences reflected a strong signature of highly restricted gene flow due to the geographical and/or climatic isolations (also see review by Qiu *et al.*³² and references therein). Indeed, our AMOVA results demonstrated moderate to high Φ_{ST} values between the two clades (ranging from $\Phi_{ST}=0.2023$ (nSSR) to $\Phi_{ST}=0.8244$ (CHS), Table 1). The cluster analysis indicated some weak admixture between the clades (Figs 3c and 4b; Supplementary Figs S5b and S7), implying stable differentiation between the two clades. However, there was some discordance in the Bayesian clustering when using different molecular markers. For instance, nSSR data clustered the populations SCPW, SCSP, SCML, SCWC, SCXL, SCKD and SCMG in the eastern clade, while nuclear sequence data (including the CHS gene), assigned these populations to the western clade. The discrepancies between nuclear sequence and nSSR data may be caused in part by their difference in mutation rates. Generally, the mean mutation rate is 10^{-9} for nuclear genes and 10^{-4} for SSR loci. This high mutation rate makes nSSR data more suitable for determining the intraspecific genetic structure^{33,34}.

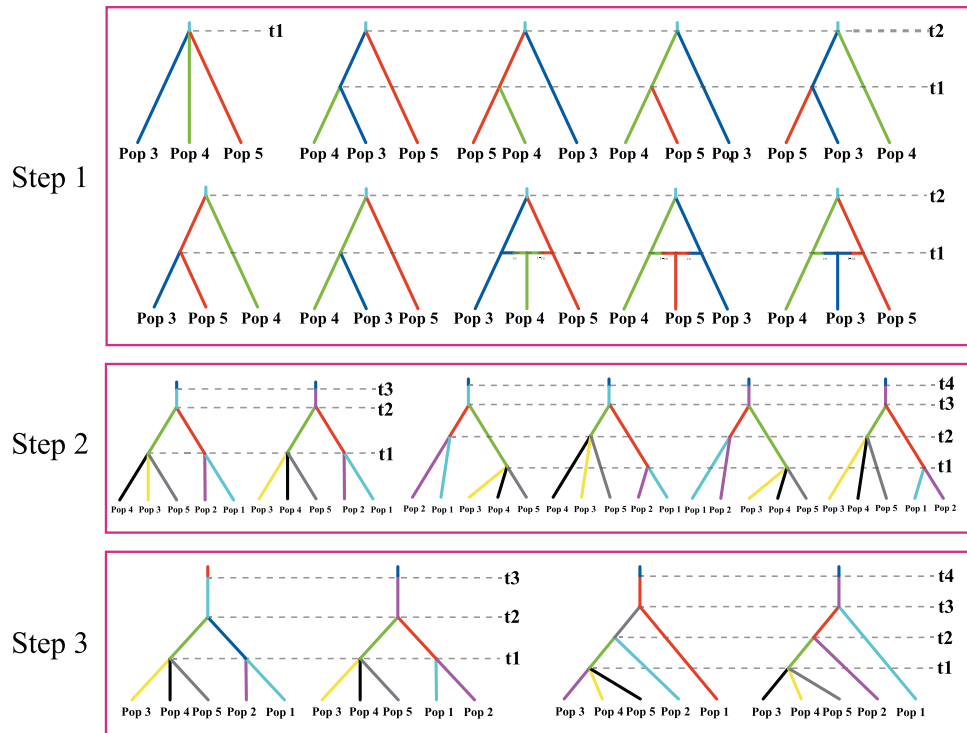


Figure 5. ABC simulation of the evolutionary history of the *Rheum palmatum* complex. Step 1: ten scenarios for the split and admixture events of three subpopulations in the eastern clade. Step 2: six scenarios for the split between and within the western and eastern clades. Step 3: four scenarios for the split among five subpopulations when considering different dispersal events.

ABC simulations suggested that the separation between the clades happened *c.* 2.0 Ma ago with subsequent divergence within clades *c.* 1.0 Ma ago. It is possible that the major uplift of the eastern edge of the QTP during the Miocene and Late Pliocene³⁵, associated with repeated climate changes in the Quaternary^{17,36}, played a critical role in the differentiation of the clades. However, the possibility that the QTP uplift alone could have contributed to the intraspecific divergence of the *R. palmatum* complex cannot be completely ruled out^{21,37–40}. For instance, the recent uplift of the QTP (*c.* 0.9–1.1 Ma)⁴¹ is proposed to have been one factor causing genetic divergence of other temperate species in this region^{32,42}. However, there are some uncertainties about the time estimates such as the mutation model and homoplasy as well as the assumption of no gene flow when using DIYABC in our study, which might underestimate divergence times over large timescales^{34,43}. Nevertheless, our time estimates in the present study are consistent with previous studies (e.g., the increased speciation rates in the Hengduan Mountains and the evolutionary radiation of *Rheum* during the Quaternary²¹), indicating that our results reflect the rate of evolutionary divergence within the *R. palmatum* complex.

The results from the ENMs demonstrated that the two lineages expanded their ranges from LIG to LGM but indicated different demography histories thereafter (e.g., western populations expanded while eastern populations shrunk, Fig. 6a), and our Bayesian skyline plots of cpDNA and nuclear genes also indicated the species complex increased its effective population size during the Quaternary (Supplementary Fig. S8). We speculate that this unusual demography of the *R. palmatum* complex might have been due to increased human activities in the Himalaya–Hengduan Mountains during recent decades which caused a decrease in suitable habitats leading to higher levels of genetic drift and more bottlenecks (Supplementary Table S5). The complicated topography in these regions provide suitable microenvironments, buffering the impacts of climatic oscillations. Thus, *R. palmatum* complex populations could have moved northwestward following their optimal conditions and expanded their ranges during the LGM. In addition, our results indicate that the two lineages occupied different areas when only areas with moderate suitability scores are considered (e.g., >0.75, Fig. 6a). This is consistent with the results from the identity test and the PCA analysis (Fig. 6b,c) showing significant differences between the two clades. Such ecological niche partitioning will have reinforced the divergence of the two clades following initial spatial separation and the adaptation of populations to their local environment, which may ultimately lead to reproductive isolation and the generation of new species if given sufficient time^{44–46}.

The present study used a chemical approach and genetic markers (cpDNA, nDNA, and nSSR) to better understand the genetic differentiation at the population level of rhubarb plants throughout its distribution range. The analysis of both chemical composition and genetic markers showed that two distinct groups are present in the *R. palmatum* complex. These two groups coincide with the areas which are considered to produce top-geoherbs and non-top-geoherbs for rhubarb in China. Ancient Chinese practitioners classed rhubarb collected from the Himalaya–Hengduan Mountains (Qinghai, Gansu, and Sichuan Provinces) as a top-geoherb, which corresponds to the populations from our western clade (e.g. QHMH, QHTJ, QHDR, GSSD, GSDB, SCXL, and SCML)

	Parameter	N1	N2	N3	N4	N5	NA	N1b	N2b	t1 ^a	t2 ^a	t3 ^a	μ	P
Scenario 1	median	3.88×10^5	6.42×10^5	6.50×10^5	6.09×10^5	6.61×10^5	4.66×10^5	6.00×10^5	—	2.41×10^5	4.14×10^5	1.47×10^6	2.84×10^{-5}	0.568
	q050	1.36×10^5	2.93×10^5	3.03×10^5	2.74×10^5	3.19×10^5	4.45×10^4	7.87×10^4	—	5.60×10^4	1.28×10^5	4.34×10^5	1.37×10^{-5}	0.246
	q950	7.93×10^5	9.25×10^5	9.22×10^5	9.16×10^5	9.27×10^5	9.44×10^5	9.62×10^5	—	1.07×10^6	1.93×10^6	2.81×10^6	8.09×10^{-5}	0.878
Scenario 2	median	5.73×10^5	4.56×10^5	6.56×10^5	6.23×10^5	6.09×10^5	3.33×10^5	—	6.10×10^5	2.13×10^5	3.68×10^5	1.31×10^6	2.89×10^{-5}	0.471
	q050	2.44×10^5	1.71×10^5	3.19×10^5	2.82×10^5	2.87×10^5	2.49×10^4	—	8.92×10^4	4.99×10^4	1.08×10^5	3.63×10^5	1.39×10^{-5}	0.168
	q950	8.99×10^5	8.32×10^5	9.27×10^5	9.20×10^5	9.09×10^5	9.05×10^5	—	9.64×10^5	9.67×10^5	1.65×10^6	2.77×10^6	8.28×10^{-5}	0.836

Table 2. Posterior median estimate and 95% highest posterior density interval (HPDI) for demographic parameters in scenarios 1 and 2 in STEP3 based on the nuclear multilocus microsatellite data for whole populations of the *Rheum palmatum* complex. ^aThe unit of timing is generation. μ: mutation rate (per generation per locus). P represents the proportion of multiple step mutations in the generalized stepwise model, GSM.

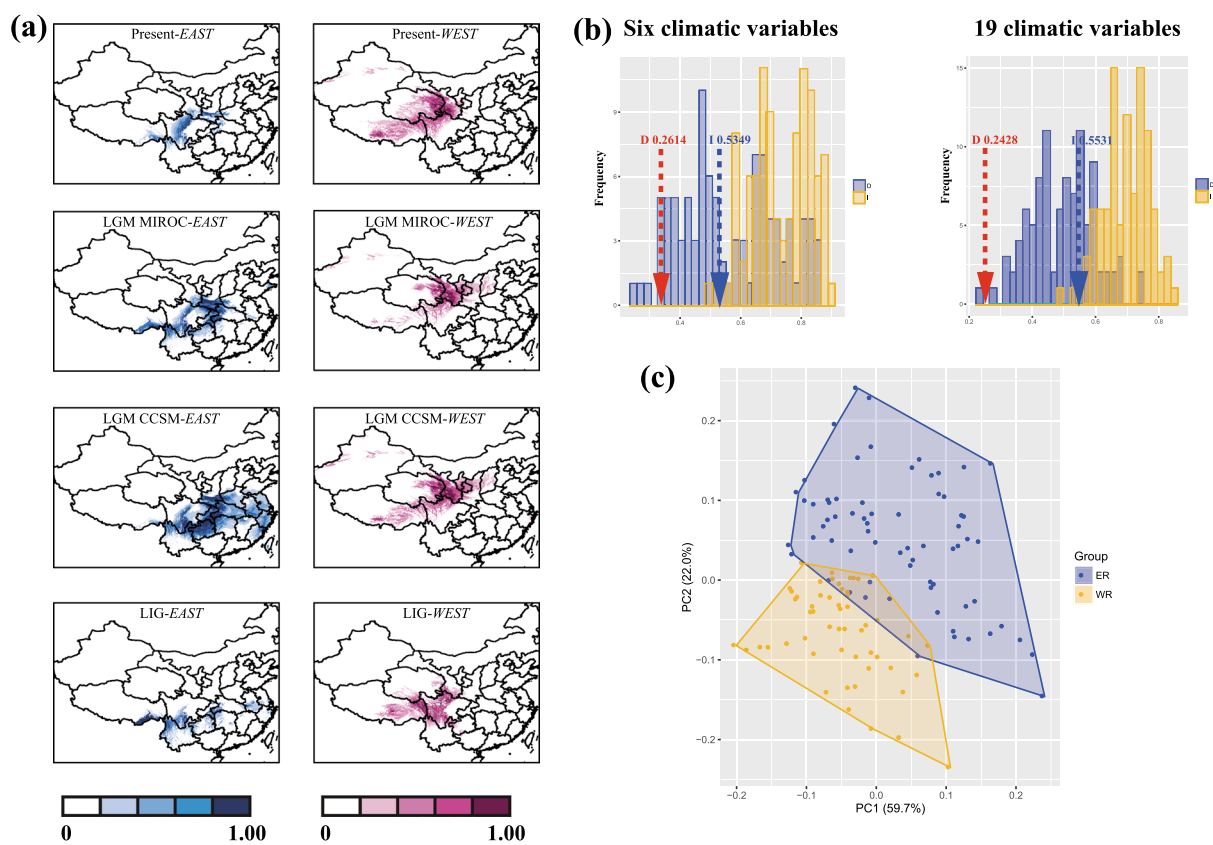


Figure 6. Results of predicted suitability areas, niche differentiation, and environmental differences in western and eastern clades of the *Rheum palmatum* complex. (a) LIG, the last interglacial (*c.* 130 Ka BP); LGM-CCSM and LGM-MIROC, last glacial maximum (*c.* 21 Ka BP) employing two different paleoclimate layers, the Community Climate System Model (CCSM) and the Model for Interdisciplinary Research on Climate (MIROC); PRE, present conditions (*c.* 1950–2000). Darker colors indicate higher probabilities of suitable climatic conditions; (b) The results of identity tests between western and eastern clades based on six environmental variables and all 19 environmental variables. Bars indicate the null distributions of *D* and *I*. Both are generated from 100 randomizations. The x-axis indicates values of *I* and *D* and the y-axis indicates number of randomizations. The arrow indicates values in actual Maxent runs; (c) Principal component analysis (PCA) plot of 20 environmental variables between the western (light blue color) and eastern (pink color) clades. The base map was drawn using ArcGIS v.10.2 (ESRI, Redlands, CA, USA).

collected in Qinghai, Gansu, and Sichuan. This is clearly shown by the cluster analyses of nSSR and nuclear genes. Our results show that chemical variation and genetic differentiation of rhubarb populations between top-geoherb and non-top-geoherb regions are consistent with the concept of geo-herbalism. Our study revealed the contribution of the genetic factors to the formation of top-geoherbs, and also showed that the formation of top-geoherbs might be closely related to species differentiation. Future studies should focus on quantitative genetics to further understand the molecular mechanism of top-geoherbs.

Overall, our results clarified that the current populations of the *R. palmatum* complex comprises two major clades (eastern and western). Quaternary environmental changes and/or the uplift of the QTP profoundly influenced the evolutionary and population demographic history of the *R. palmatum* complex as well as its current genetic structure. Chemical and genetic variation exists within the *R. palmatum* complex, and are consistent with the top-geoherb and non-top-geoherb areas of rhubarb. These findings reveal that genetic differentiation is at the core of the geo-herbalism concept.

Materials and Methods

Sample collection and DNA extraction. Leaf samples were collected from 38 populations covering the whole range of the *R. palmatum* complex (Supplementary Table S7 and Fig. S3a). All samples were assayed for 14 nSSR loci ($n = 479$), and a subset of samples were sequenced for cpDNA and nDNA regions ($n = 377$). In addition, the underground roots and rhizomes of thirty-one populations of the *R. palmatum* complex were collected to minimize the effect of seasonal changes on the concentrations of chemical constituents in the wild⁴⁷ and their chemical compounds were analyzed using three individuals to reduce individual differences.

Total DNA was isolated from dried leaf tissue using a plant total genomic DNA kit (Tiangen, Beijing, China). The quality and quantity of DNA was determined on 1% TAE agarose gels and with a NanoDrop® ND-1000 spectrophotometer (Thermo Fisher Scientific, Wilmington, DE, USA), respectively. The DNA was diluted to a final concentration of approximately 20 ng/ μ L.

Analysis of chemical compounds. Dried root and rhizome samples were pulverized to extract free anthraquinones and anthraquinone glycosides. The powder was filtered through 180 μ m sieves and accurately weighed when fine (800 mg). Then, 70% methanol (25 mL) was added and the mixture was weighed again. The mixture was extracted by ultrasonication for 1 h (40 °C, 300 w). After cooling, 70% methanol was added and transferred to a 25 mL volumetric flask. The supernatant fluid was filtered through a 0.22 μ m membrane filter, and the filtrate was analyzed using high-performance liquid chromatography (HPLC) (see details in Note S1).

Chloroplast and nuclear DNA sequencing and microsatellite genotyping. A total of six cpDNA regions (*trnL-trnF*, *ndhJ-trnE*, *psaI-accD*, *rpl20-rps12*, *rpl32-trnL*, and *psbA-trnH*) and five nuclear genes (*CHS* and the four genes R27, R30, R46, R56 which were developed from transcriptome data) were sequenced. Primers and PCR amplification of the cpDNA regions were described in previous studies^{48,49} (Supplementary Table S8), and the amplification of five nDNA regions was described in Ma *et al.*⁵⁰. Sequences were generated (excepted for the *CHS* gene) using an ABI 377 DNA sequencer (Sangon Biological Engineering Technology & Service Co.). Clone sequencing was used for the *CHS* gene. More than ten colonies of each sample were randomly chosen and sequenced to survey sequence variations in multiple copies and then the program MEGA7⁵¹ was employed to align and edit sequences. All newly acquired sequences have been submitted to GenBank with accession numbers ranged from MH457640 to MH457707 and MH465255 to MH465390.

All sampled individuals were genotyped at 14 nSSR loci (see Supplementary Table S3) developed from the genomic resources of *R. palmatum*. The PCR conditions of nSSR loci were similar to the cpDNA fragments, except for the annealing temperature. The PCR products were genotyped using a 3730XL automated Genetic Analyzer (Applied Biosystems, Foster City, CA). Allele sizes were determined in GENEMAPPER (version 3.7, Applied Biosystems).

Sequence analysis. The haplotype (H_d) and nucleotide (π) diversity, the number of segregating sites (S), the Watterson's parameter (θ_w)⁵², and the minimum number of recombinant events (R_m)⁵³ were estimated, and tests were carried out for departure from the neutral model based on Tajima's D ⁵⁴, and Fu and Li's D^* and F^{*55} using DNASP v5.1⁵⁶. The total gene diversity (H_T), gene diversity within populations (H_S), and coefficients of differentiation G_{ST} and N_{ST} were calculated using SPAGEDI⁵⁷ based on 10,000 random permutations. Genealogical relationships of the haplotypes were identified using POPART based on a TCS network algorithm⁵⁸. A mismatch analysis was conducted using ARLEQUIN v3.5⁵⁹ to test for spatial expansion of test populations. The goodness-of-fit was tested with the sum of squared deviations (SSD) between observed and expected mismatch distributions and Harpending's raggedness index (H_{Rag})⁶⁰. We also employed jMODELTEST v1.0⁶¹ to evaluate the best-fit model of nucleotide substitution for maximum likelihood (ML) method to infer cpDNA haplotype relationships (GTR + I + G in the present case), and used PAUP* v4.0 beta 10⁶² to determine the relationships of cpDNA haplotypes and visualized results using FIGTREE v1.3.1 (<http://tree.bio.ed.ac.uk/software/figtree/>). The program STRUCTURE v 2.3.3⁶³ was used to infer the population genetic structure of nDNA sequences based on the admixture model with allele frequencies correlated. The number of clusters (K) was set to vary from one to 20. For each value of K , we performed 20 runs with a burn-in length of 200,000 and a run length of 500,000 Markov chain Monte Carlo (MCMC) replications. The best K values were determined using the methods described by Pritchard *et al.*⁶³ and Evanno *et al.*⁶⁴. The web-based program *Structure Harvester*⁶⁵ was used for visualizing the STRUCTURE output. The estimated admixture coefficients (Q matrix) over the 20 replicates in each K were averaged using CLUMPP v1.1⁶⁶, and the graphics of optimal K were produced using DISTRUCT v1.1⁶⁷. In addition, in order to quantify variation in nSSRs and DNA sequences between and within the two main geographical regions, we performed analyses of molecular variance (AMOVA) in ARLEQUIN using 10,000 permutations. Additionally, we used BSP in BEAST v.1.7.5⁶⁸ based on the combined cpDNA data and four nuclear genes each to estimate effective population size changes within the species complex. Linear and stepwise models were explored using an uncorrelated lognormal relaxed clock. Runs consisted of 20,000,000 generations, with trees sampled every 5,000 generations. The BSP was visualized in the program Tracer v1.5, which summarizes the posterior distribution of population size over time. The cpDNA substitution rates for most angiosperm species have been estimated to vary between 1 and 3×10^{-9} substitutions per site per year (s/s/yr)⁶⁹, while those for nDNA in shrubs and herbal

plants vary between 3.46 and 8.69×10^{-9} (s/s/yr)³³. Given the uncertainties in these rate values, we used normal distribution priors with a mean of 2×10^{-9} and a SD of 6.080×10^{-10} for cpDNA, and a mean of 6.075×10^{-9} and a SD of 1.590×10^{-9} for nuclear gene to cover these rate ranges within the 95% range of the distribution for our BSP analyses.

Microsatellite data analysis. Population genetic analysis. For the nuclear microsatellite dataset, all 14 loci were checked for possible null alleles using MICRO-CHECKER v2.3.3⁷⁰. Tests for departures from Hardy–Weinberg equilibrium (HWE) and linkage disequilibrium (LD) were performed in each population and a globally unified population using FSTAT v2.9.3⁷¹ based on 1,000 permutations ($\alpha = 0.05$), and corrected for multiple tests using the sequential Bonferroni method⁷². In addition, we conducted a modified version of F_{ST} outlier loci analysis⁷³ under the assumption of the IAM using LOSITAN⁷⁴ with 50,000 simulations to generate 95% confidence intervals (CIs), and the results indicated all 14 loci were neutral, so all loci were retained for use in subsequent analyses.

For each locus, genetic diversity was assessed by calculating the observed number of alleles (A_O), the observed heterozygosity (H_O), the expected heterozygosity over all populations (H_E), the genetic diversity within the population (H_S), the total genetic diversity (H_T), and the standardized genetic differentiation G'_{ST} ⁷⁵ across loci with FSTAT and/or GENALEX v6.5⁷⁶ based on 1,000 bootstrap permutations. Compared with the traditional measures (e.g., F_{ST} and G_{ST}), G'_{ST} is a more suitable measure of differentiation for highly polymorphic markers such as microsatellites^{75,77}. For each population, we also used the program GENALEX to estimate the values of H_O , H_E , private allelic richness (P_{AR}), fixation index (F_{IS}), and allele richness (A_S) across all loci.

Using the nSSR dataset, genetic clusters were identified using a Bayesian analysis in STRUCTURE with the same parameter settings as above, except for the K was set to vary from one to 38. To test how the geographic distance affected the genetic composition of the *R. palmatum* complex, Mantel tests with 10,000 random permutations were performed between the matrix of pairwise $F_{ST}/(1 - F_{ST})$ and that of the natural logarithm of the geographic distances at species and each gene pool revealed by STRUCTURE analysis, respectively, using the package Vegan⁷⁸ in R 3.3.0⁷⁹.

Testing for genetic bottleneck and estimates for historical and contemporary gene flow. The Wilcoxon's signed rank test and the mode-shift test implemented in BOTTLENECK v1.2.02⁸⁰ was used to detect whether a population had experienced a size decrease over extended vs. more recent time scales ($\sim 2Ne-4Ne$ generations in the past vs. a few dozen generations ago⁸¹, See details in Note S1).

To estimate historical gene flow (much longer period of time, $\sim 4Ne$ generations in the past⁸²), between the STRUCTURE clusters and between sub-clusters within each STRUCTURE cluster (see results section), we used the program MIGRATE v3.6⁸³ to estimate the effective number of migrants ($2Nm$, where Ne is the effective population size and m is the migration rate per generation). Also, we estimated the contemporary gene flow using BAYESASS v3.03⁸⁴ (See details in Note S1).

Tests of demographic history by ABC modeling. In order to decrease the computational amount and time, we narrowed scenarios by defining nested subsets of competing scenarios that were analyzed sequentially based on STRUCTURE analysis using the ABC procedure⁸⁵ as performed in DIYABC v2.1.0^{86,87}, because hundreds of scenarios can be tested with five populations. For the ABC analysis, ten alternative scenarios of population history for the three sub-clusters composing the eastern clade were summarized (step 1 in Fig. 1; Supplementary Table S9) and tested using DIYABC, and the population sizes for all sub-clusters (pops1-5) were set to be constant in the analysis excepted for the ancestral population size. Then, we analyzed all five sub-clusters, and six alternative scenarios were included in this step (step 2 in Fig. 1; Supplementary Table S9), considering the best scenario in step 1 (i.e., scenario 1). Finally, given that the *R. palmatum* complex originated on the QTP¹⁵, we set four competitive scenarios (including the best scenarios in step 2, i.e., scenarios 1 and 2) to determine its migration route(s), and whether the genetic disconnection between the two clades was old or recent (step 3 in Fig. 1; Supplementary Table S9). A reference table was generated containing 20,000,000 simulated data sets for all scenarios (on average 1,000,000 per scenario), assuming uniform priors on all parameters for each scenario (Supplementary Table S10). A goodness-of-fit test was used to check the priors of all parameters before implementing the simulations. Each simulation was summarized by the following summary statistics: the mean number of alleles and the mean genetic diversity for each clade, and the F_{ST} , the mean classification index, and shared allele distance between pairs of clades. After comparing the posterior probability of scenarios in different steps using the logistic regression and direct approaches, the posterior distribution of historical demographic parameters for the final step was estimated using the 1% simulated datasets closest to the observed dataset for the local linear regression. The average generation time was assumed to be five years according to our field observations for the *R. palmatum* complex.

Ecological niche models (ENMs) and environmental factor analysis. We used ENMs to infer suitable climate envelopes for the two clades of the *Rheum* species complex (see result section) in the present, the Last Glacial Maximum (LGM: ca. 21,000 years before present; BP), and the last interglacial (LIG: ca. 120,000–140,000 years BP) using the maximum entropy method implemented in MAXENT v3.3.3^{88,89}, and used ENMTools v1.3^{90,91} to calculate Schoener's D ⁹² and standardized Hellinger distance (calculated as I) to measure the niche similarity between clades (See details in Note S1).

In order to evaluate the effect of present climatic conditions on the genetic differentiation between the main two geographical regions, three methods were used to evaluate the impacts of climate on population differentiation. First, the lineage-level divergence associated with each of the environmental variables was examined using one-way nonparametric ANOVA (i.e., Kruskal–Wallis test⁹³) in R, and the distributions of all 20 environmental

factors each in the two clades were displayed in kernel density plots using the R package ggplot⁹⁴. Then also using R, we performed a permutational MANOVA analysis for all environmental variables simultaneously for the two clades to evaluate whether environmental conditions differed significantly between their sites of occurrence. Finally, we employed the R package ggfortify⁹⁵ to perform a principal component analysis (PCA) to determine whether the genetic groups were ecologically differentiated.

References

- Hu, S. *Geoherbs in China*. (Heilongjiang Science and Technology Press, 1989).
- Xie, Z. Discussion about geoherbs. *J. Tradit. Chin. Med.* **40**, 43–46 (1990).
- Huang, L., Guo, L., Ma, C., Gao, W. & Yuan, Q. Top-geoherbs of traditional Chinese medicine: common traits, quality characteristics and formation. *Front. Med.* **5**, 185–194 (2011).
- Zhao, Z., Guo, P. & Brand, E. The formation of daodi medicinal materials. *J. Ethnopharmacol.* **140**, 476–481 (2012).
- Huang, L., Xiao, P., Guo, L. & Gao, W. Molecular pharmacognosy. *Sci. China Life Sci.* **53**, 643–652 (2010).
- Pan, F. Daodi medicinal material is the essence of Chinese medicine—a review of the 390th session of Xiangshan Science Conference. Science Times. Paper read at Science Times, at Beijing, China. (2011).
- Nicotra, A. B. *et al.* Population and phylogenomic decomposition via genotyping-by-sequencing in Australian *Pelargonium*. *Mol. Ecol.* **25**, 2000–2014 (2016).
- Fujita, M. K., Leaché, A. D., Burbrink, F. T., McGuire, J. A. & Moritz, C. Coalescent-based species delimitation in an integrative taxonomy. *Trends Ecol. Evol.* **27**, 480–488 (2012).
- Ross, K. G., Gotzek, D., Ascunce, M. S. & Shoemaker, D. D. Species delimitation: a case study in a problematic ant taxon. *Syst. Biol.* **59**, 162–184 (2009).
- Renema, W. *et al.* Hopping hotspots: global shifts in marine biodiversity. *science* **321**, 654–657 (2008).
- Thompson, J.N. *The geographic mosaic of coevolution*. (University of Chicago Press, 2005).
- Zhang, Z. *et al.* Comparative study on quality of Flos Lonicerae between geo-authentic and non-authentic producing areas. *China journal of Chinese materia medica* **32**, 786–788 (2007).
- Chen, W., Liu, Y., Qiao, C. & Liu, Y. The chemical compounds and the content of sarsasapogenin of *Anemarrhena asphodeloides* Bunge. from different producing areas. *Acad. J. Second Military Med. Univ.* **18**, 84–86 (1997).
- Guo, L.-P., Huang, L.-Q., Jiang, Y.-X. & Zhan, Y.-H. RAPD analysis on genetic structure of *Atractylodes lancea*. *Chinese Pharmaceutical Journal* **3**, 171–181 (2006).
- Yuan, Q.-J. *et al.* Impacts of recent cultivation on genetic diversity pattern of a medicinal plant, *Scutellaria baicalensis* (Lamiaceae). *BMC Genet.* **11**, 29 (2010).
- Davis, M. B. & Shaw, R. G. Range shifts and adaptive responses to Quaternary climate change. *Science* **292**, 673–679 (2001).
- Hewitt, G. M. Genetic consequences of climatic oscillations in the Quaternary. *Philos. T. R. Soc. B* **359**, 183–195 (2004).
- Jia, D.-R. *et al.* Out of the Qinghai–Tibet Plateau: evidence for the origin and dispersal of Eurasian temperate plants from a phylogeographic study of *Hippophaë rhamnoides* (Elaeagnaceae). *New Phytol.* **194**, 1123–1133 (2012).
- Sun, H. *et al.* Survival and reproduction of plant species in the Qinghai–Tibet Plateau. *J. Syst. Evol.* **52**, 378–396 (2014).
- Wang, H. *et al.* Phylogeographic structure of *Hippophae tibetana* (Elaeagnaceae) highlights the highest microrefugia and the rapid uplift of the Qinghai–Tibetan Plateau. *Mol. Ecol.* **19**, 2964–2979 (2010).
- Sun, Y., Wang, A., Wan, D., Wang, Q. & Liu, J. Rapid radiation of *Rheum* (Polygonaceae) and parallel evolution of morphological traits. *Mol. Phylogenet. Evol.* **63**, 150–158 (2012).
- Bao, B. & Grabovskaya-Borodina, A. *Rheum*. (Beijing and St. Louis: Science Press and Missouri Botanical Garden, 2003).
- Wang, X.-M. *et al.* Genetic diversity of the endemic and medicinally important plant *Rheum officinale* as revealed by inter-simple sequence repeat (ISSR) markers. *Int. J. Mol. Sci.* **13**, 3900–3915 (2012).
- Wang, X.-M. *et al.* Genetic variation in *Rheum palmatum* and *Rheum tanguticum* (Polygonaceae), two medicinally and endemic species in China using ISSR markers. *PLoS One* **7**, e11667 (2012).
- Kimura, M. K. *et al.* Evidence for cryptic northern refugia in the last glacial period in *Cryptomeria japonica*. *Ann. Bot.-London.* **114**, 1687–1700 (2014).
- Li, L. *et al.* Pliocene intraspecific divergence and Plio-Pleistocene range expansions within *Picea likiangensis* (Lijiang spruce), a dominant forest tree of the Qinghai–Tibet Plateau. *Mol. Ecol.* **22**, 5237–5255 (2013).
- Mayol, M. *et al.* Adapting through glacial cycles: insights from a long-lived tree (*Taxus baccata*). *New Phytol.* **208**, 973–986 (2015).
- Wang, X.-M., Hou, X.-Q., Zhang, Y.-Q. & Li, Y. Morphological variation in leaf dissection of *Rheum palmatum* complex (Polygonaceae). *PLoS One* **9**, e110760 (2014).
- Wang, X.-M., Hou, X.-Q., Zhang, Y.-Q. & Li, Y. Distribution pattern of genuine species of rhubarb as traditional Chinese medicine. *J. Med. Plants Res.* **4**, 1865–1876 (2010).
- Wang, J.-F., Gong, X., Chiang, Y.-C. & Kuroda, C. Phylogenetic patterns and disjunct distribution in *Ligularia hodgsonii* Hook. (Asteraceae). *J. Biogeogr.* **40**, 1741–1754 (2013).
- Sun, Y.-X. *et al.* Chloroplast phylogeography of the East Asian Arcto-Tertiary relict *Tetracentron sinense* (Trochodendraceae). *J. Biogeogr.* **41**, 1721–1732 (2014).
- Qiu, Y.-X., Fu, C.-X. & Comes, H. P. Plant molecular phylogeography in China and adjacent regions: Tracing the genetic imprints of Quaternary climate and environmental change in the world's most diverse temperate flora. *Mol. Phylogenet. Evol.* **59**, 225–244 (2011).
- Richardson, J. E., Pennington, R. T., Pennington, T. D. & Hollingsworth, P. M. Rapid diversification of a species-rich genus of neotropical rain forest trees. *Science* **293**, 2242–2245 (2001).
- Selkoe, K. A. & Toonen, R. J. Microsatellites for ecologists: a practical guide to using and evaluating microsatellite markers. *Ecol. Lett.* **9**, 615–629 (2006).
- Sun, B.-N. *et al.* Reconstructing Neogene vegetation and climates to infer tectonic uplift in western Yunnan, China. *Palaeogeogr. Palaeoclimatol.* **304**, 328–336 (2011).
- Hewitt, G. The genetic legacy of the Quaternary ice ages. *Nature* **405**, 907–913 (2000).
- Li, J.-J. & Fang, X.-M. Uplift of the Tibetan Plateau and environmental changes. *Chinese Sci. Bull.* **44**, 2117–2124 (1999).
- Renner, S. S. Available data point to a 4-km-high Tibetan Plateau by 40 Ma, but 100 molecular-clock papers have linked supposed recent uplift to young node ages. *J. Biogeogr.* **43**, 1479–1487 (2016).
- Wang, C. *et al.* Constraints on the early uplift history of the Tibetan Plateau. *P. Natl. Acad. Sci. USA* **105**, 4987–4992 (2008).
- Royden, L. H., Burchfiel, B. C. & van der Hilst, R. D. The geological evolution of the Tibetan Plateau. *Science* **321**, 1054–1058 (2008).
- Sun, J. & Liu, T. Stratigraphic evidence for the uplift of the Tibetan Plateau between ~1.1 and ~0.9 myr ago. *Quaternary Res.* **54**, 309–320 (2000).
- Wen, J., Zhang, J., Nie, Z.-L., Zhong, Y. & Sun, H. Evolutionary diversifications of plants on the Qinghai–Tibetan Plateau. *Front. Genet.* **5** (2014).
- Estoup, A., Jarne, P. & Cornuet, J.-M. Homoplasy and mutation model at microsatellite loci and their consequences for population genetics analysis. *Mol. Ecol.* **11**, 1591–1604 (2002).

44. Liu, J. *et al.* Geological and ecological factors drive cryptic speciation of yews in a biodiversity hotspot. *New Phytol.* **199**, 1093–1108 (2013).
45. Nosil, P., Harmon, L. J. & Seehausen, O. Ecological explanations for (incomplete) speciation. *Trends Ecol. Evol.* **24**, 145–156 (2009).
46. Rundle, H. D. & Nosil, P. Ecological speciation. *Ecol. Lett.* **8**, 336–352 (2005).
47. Paneitz, A. & Westendorf, J. Anthranoid contents of rhubarb (*Rheum undulatum* L.) and other *Rheum* species and their toxicological relevance. *Eur. Food Res. Technol.* **210**, 97–101 (1999).
48. Taberlet, P., Gielly, L., Pautou, G. & Bouvet, J. Universal primers for amplification of three non-coding regions of chloroplast DNA. *Plant Mol. Biol.* **17**, 1105–1109 (1991).
49. Hamilton, M. Four primer pairs for the amplification of chloroplast intergenic regions with intraspecific variation. *Mol. Ecol.* **8**, 521–523 (1999).
50. Ma, L.-Q. *et al.* A novel type III polyketide synthase encoded by a three-intron gene from *Polygonum cuspidatum*. *Planta* **229**, 457 (2009).
51. Kumar, S., Stecher, G. & Tamura, K. MEGA7: Molecular Evolutionary Genetics Analysis Version 7.0 for Bigger Datasets. *Mol. Biol. Evol.* **33**, 1870–1874 (2016).
52. Watterson, G. On the number of segregating sites in genetical models without recombination. *Theor. Popul. Biol.* **7**, 256–276 (1975).
53. Hudson, R. R. & Kaplan, N. L. Statistical properties of the number of recombination events in the history of a sample of DNA sequences. *Genetics* **111**, 147–164 (1985).
54. Tajima, F. Statistical method for testing the neutral mutation hypothesis by DNA polymorphism. *Genetics* **123**, 585–595 (1989).
55. Fu, Y. X. & Li, W. H. Statistical tests of neutrality of mutations. *Genetics* **133**, 693–709 (1993).
56. Librado, P. & Rozas, J. DnaSPv5: a software for comprehensive analysis of DNA polymorphism data. *Bioinformatics* **25**, 1451–1452 (2009).
57. Hardy, O. J. & Vekemans, X. SPAGeDi: a versatile computer program to analyse spatial genetic structure at the individual or population levels. *Mol. Ecol. Notes* **2**, 618–620 (2002).
58. Leigh, J. W. & Bryant, D. POPART: full-feature software for haplotype network construction. *Methods Ecol. Evol.* **6**, 1110–1116 (2015).
59. Excoffier, L. & Lischer, H. E. L. Arlequin suite ver 3.5: a new series of programs to perform population genetics analyses under Linux and Windows. *Mol. Ecol. Resour.* **10**, 564–567 (2010).
60. Harpending, H. C. Signature of ancient population growth in a low-resolution mitochondrial DNA mismatch distribution. *Hum. Biol.* **66**, 591–600 (1994).
61. Posada, D. jModelTest: Phylogenetic model averaging. *Mol. Biol. Evol.* **25**, 1253–1256 (2008).
62. Swofford, D. L. PAUP*: phylogenetic analysis using parsimony, version 4.0 b10. (2003).
63. Pritchard, J. K., Stephens, M. & Donnelly, P. Inference of population structure using multilocus genotype data. *Genetics* **155**, 945–959 (2000).
64. Evanno, G., Regnaut, S. & Goudet, J. Detecting the number of clusters of individuals using the software structure: a simulation study. *Mol. Ecol.* **14**, 2611–2620 (2005).
65. Earl, D. & vonHoldt, B. Structure Harvester: a website and program for visualizing Structure output and implementing the Evanno method. *Conserv. Genet. Resour.* **4**, 359–361 (2012).
66. Jakobsson, M. & Rosenberg, N. A. Clumpp: a cluster matching and permutation program for dealing with label switching and multimodality in analysis of population structure. *Bioinformatics* **23**, 1801–1806 (2007).
67. Rosenberg, N. A. Distruct: a program for the graphical display of population structure. *Mol. Ecol. Notes* **4**, 137–138 (2004).
68. Drummond, A. J. & Rambaut, A. Beast: Bayesian evolutionary analysis by sampling trees. *BMC Evol. Biol.* **7**, 214 (2007).
69. Wolfe, K. H., Li, W.-H. & Sharp, P. M. Rates of nucleotide substitution vary greatly among plant mitochondrial, chloroplast, and nuclear DNAs. *P. Natl. Acad. Sci. USA* **84**, 9054–9058 (1987).
70. Van Oosterhout, C., Hutchinson, W. F., Wills, D. P. & Shipley, P. Micro-Checker: software for identifying and correcting genotyping errors in microsatellite data. *Mol. Ecol. Notes* **4**, 535–538 (2004).
71. Goudet, J. Fstat; a program to estimate and test gene diversities and fixation indices version 2.9.3. <http://www.unil.ch/izea/software/fstat.html> (2001).
72. Rice, W. R. Analyzing tables of statistical tests. *Evolution* **223**–225 (1989).
73. Beaumont, M. A. & Nichols, R. A. Evaluating Loci for Use in the Genetic Analysis of Population Structure. *P. Roy. Soc. B-Biol. Sci.* **263**, 1619 (1996).
74. Antao, T., Lopes, A., Lopes, R., Beja-Pereira, A. & Luikart, G. Lositan: A workbench to detect molecular adaptation based on a *Fst*-outlier method. *BMC Bioinformatics* **9**, 323 (2008).
75. Hedrick, P. W. A standardized genetic differentiation measure. *Evolution* **59**, 1633–1638 (2005).
76. Peakall, R. & Smouse, P. E. GenALEX 6.5: genetic analysis in Excel. Population genetic software for teaching and research—an update. *Bioinformatics* **28**, 2537–2539 (2012).
77. Crawford, N. G. Smogd: software for the measurement of genetic diversity. *Mol. Ecol. Resour.* **10**, 556–557 (2010).
78. Oksanen, J. *et al.* The vegan package. *Community ecology package* **10**, 631–637 (2007).
79. Team, R. C. R: A language and environment for statistical computing [Internet]. Vienna, Austria: R Foundation for Statistical Computing; 2013. *Document freely available on the internet at:* <http://www.r-project.org> (2015).
80. Piry, S., Luikart, G. & Cornuet, J. M. Bottleneck: a computer program for detecting recent reductions in the effective size using allele frequency data. *J. Hered.* **90**, 502–503 (1999).
81. Luikart, G. & Cornuet, J.-M. Empirical evaluation of a test for identifying recently bottlenecked populations from allele frequency data. *Conserv. Biol.* **12**, 228–237 (1998).
82. Beerli, P. & Felsenstein, J. Maximum likelihood estimation of a migration matrix and effective population sizes in *n* subpopulations by using a coalescent approach. *P. Natl. Acad. Sci. USA* **98**, 4563–4568 (2001).
83. Beerli, P. Comparison of Bayesian and maximum-likelihood inference of population genetic parameters. *Bioinformatics* **22**, 341–345 (2006).
84. Wilson, G. A. & Rannala, B. Bayesian inference of recent migration rates using multilocus genotypes. *Genetics* **163**, 1177–1191 (2003).
85. Beaumont, M. A., Zhang, W. & Balding, D. J. Approximate Bayesian computation in population genetics. *Genetics* **162**, 2025–2035 (2002).
86. Cornuet, J.-M. *et al.* Inferring population history with DIYABC: a user-friendly approach to approximate Bayesian computation. *Bioinformatics* **24**, 2713–2719 (2008).
87. Cornuet, J.-M. *et al.* DIYABCv2.0: a software to make approximate Bayesian computation inferences about population history using single nucleotide polymorphism, DNA sequence and microsatellite data. *Bioinformatics* (2014).
88. Phillips, S. J. & Dudik, M. Modeling of species distributions with Maxent: new extensions and a comprehensive evaluation. *Ecography* **31**, 161–175 (2008).
89. Phillips, S. J., Anderson, R. P. & Schapire, R. E. Maximum entropy modeling of species geographic distributions. *Ecol. Model.* **190**, 231–259 (2006).
90. Warren, D. L., Glor, R. E. & Turelli, M. Environmental niche equivalency versus conservatism: quantitative approaches to niche evolution. *Evolution* **62**, 2868–2883 (2008).

91. Warren, D. L., Glor, R. E. & Turelli, M. ENMTools: a toolbox for comparative studies of environmental niche models. *Ecography* **33**, 607–611 (2010).
92. Schoener, T. W. The Anolis lizards of Bimini: resource partitioning in a complex fauna. *Ecology*, 704–726 (1968).
93. Conover, W. J. *Practical nonparametric statistics*. (John Wiley & Sons, 1999).
94. Wickham, H. *ggplot2: elegant graphics for data analysis*. (Springer Science & Business Media, 2009).
95. Horikoshi, M. & Tang, Y. ggfortify: data visualization tools for statistical analysis results. R package version 0.0.4. (2015).

Acknowledgements

This work was supported by the National Natural Science Foundation of China (31470401, 31770364 and 81001602) and the China Postdoctoral Science Foundation (2018M633490).

Author Contributions

X.W. conceived and designed the experiments. X.H. and K.F. collected the materials. X.H., K.F., X.S., M.H., Y.Z., and J.S. performed the experiments. X.W., L.F., T.Z., and L.H. compiled and interpreted the data. X.W., L.F., T.Z., and M.R. wrote the manuscript. All authors reviewed and discussed the manuscript.

Additional Information

Supplementary information accompanies this paper at <https://doi.org/10.1038/s41598-018-27510-1>.

Competing Interests: The authors declare no competing interests.

Publisher's note: Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.



Open Access This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons license, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons license, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons license and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this license, visit <http://creativecommons.org/licenses/by/4.0/>.

© The Author(s) 2018