

Manica Balant*, Daniel Vitales*, Zhiqiang Wang*, Zoltán Barina, Lin Fu, Tiangang Gao, Teresa Garnatje, Airy Gras, Muhammad Qasim Hayat, Marine Oganessian, Jaume Pellicer, Seyed A. Salami, Alexey P. Seregin, Nina Stepanyan-Gandilyan, Nusrat Sultana, Shagdar Tsooj, Magsar Urgamal, Joan Vallès, Robin van Velzen, Lisa Pokorný. 2025. “**Integrating target capture with whole genome sequencing of recent and natural history collections to explain the phylogeography of wild-growing and cultivated *Cannabis***”. *Plants People Planet.* 1-18. <https://doi.org/10.1002/ppp3.70043> [* = equal contributions]

Class focus area:
Phylogeography

EEB603: Brian O'Meara

All quotes and images from the above paper unless otherwise noted

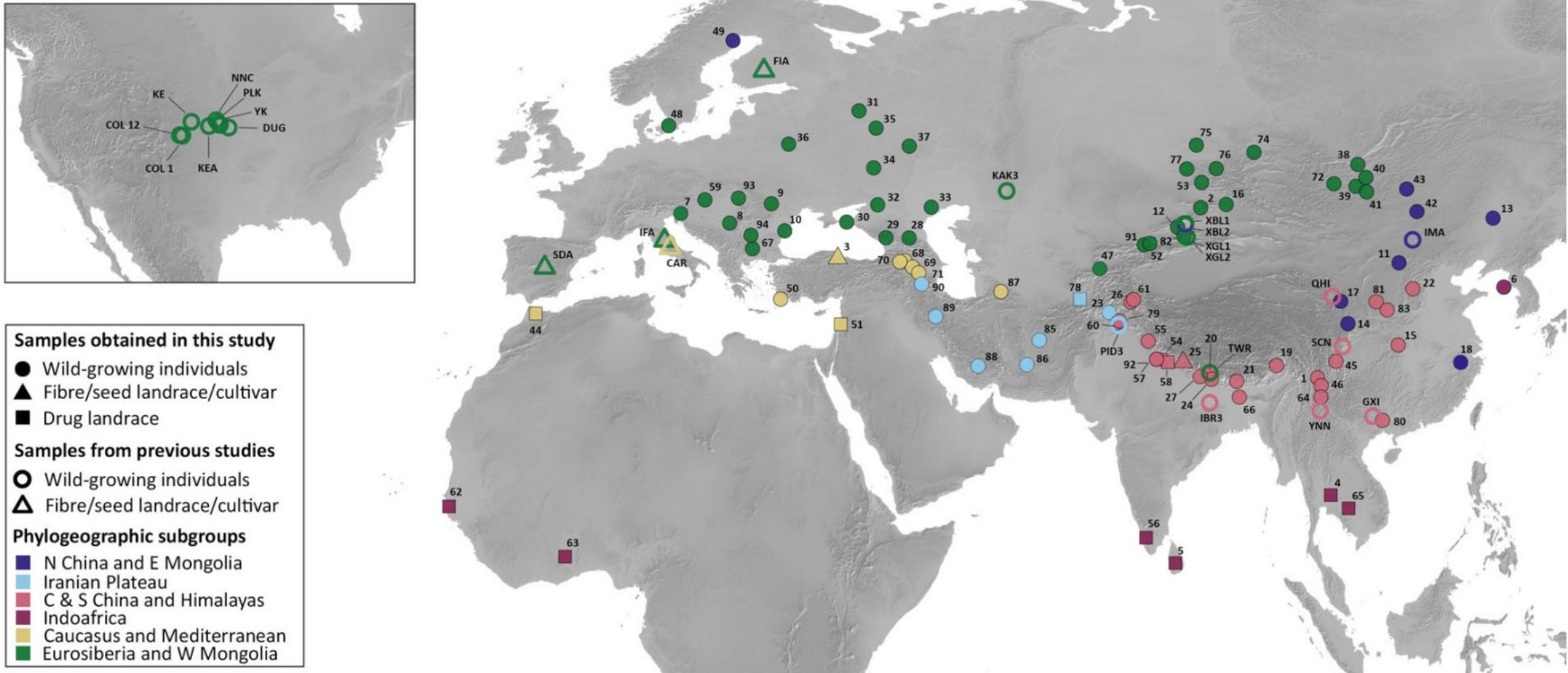
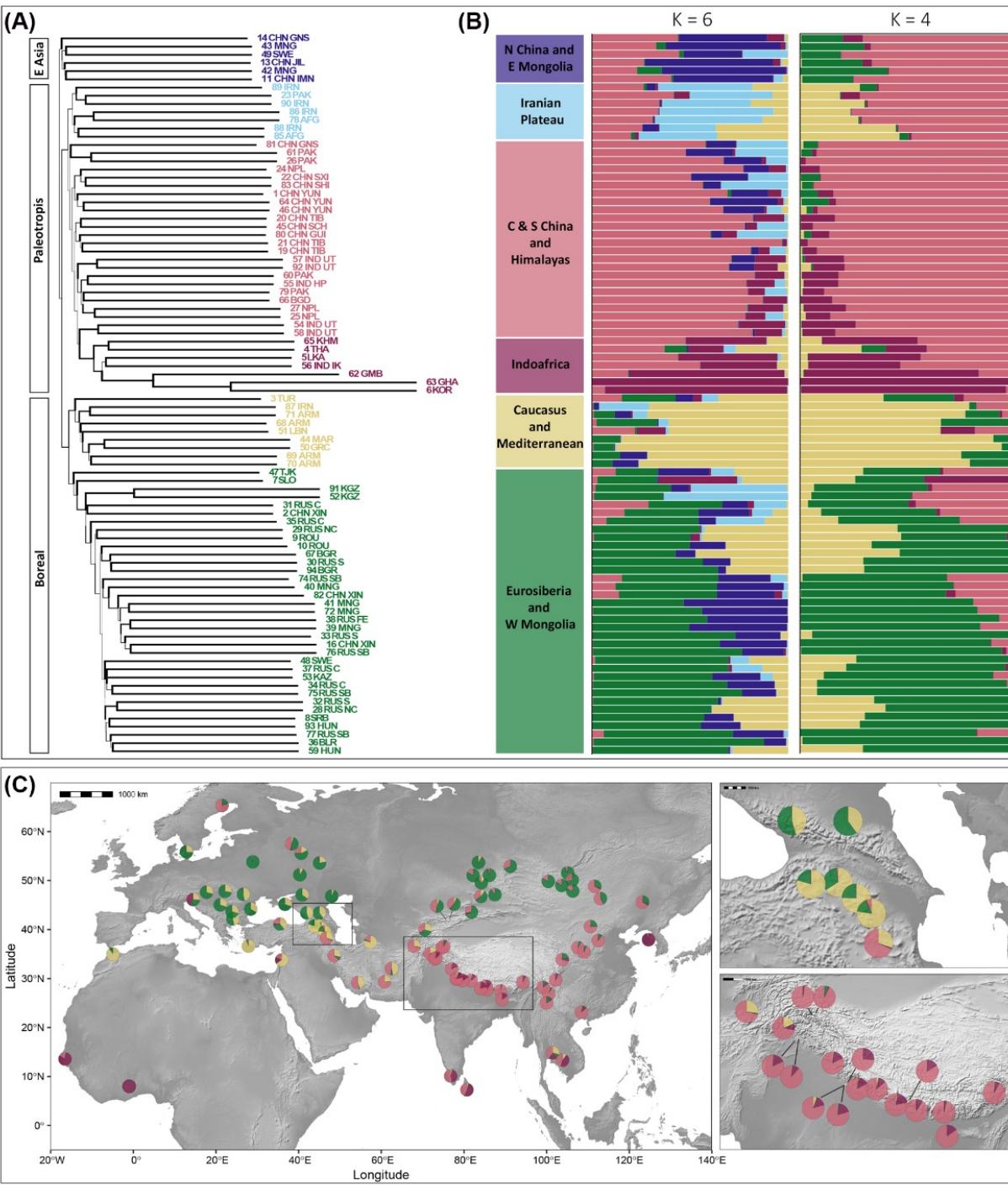
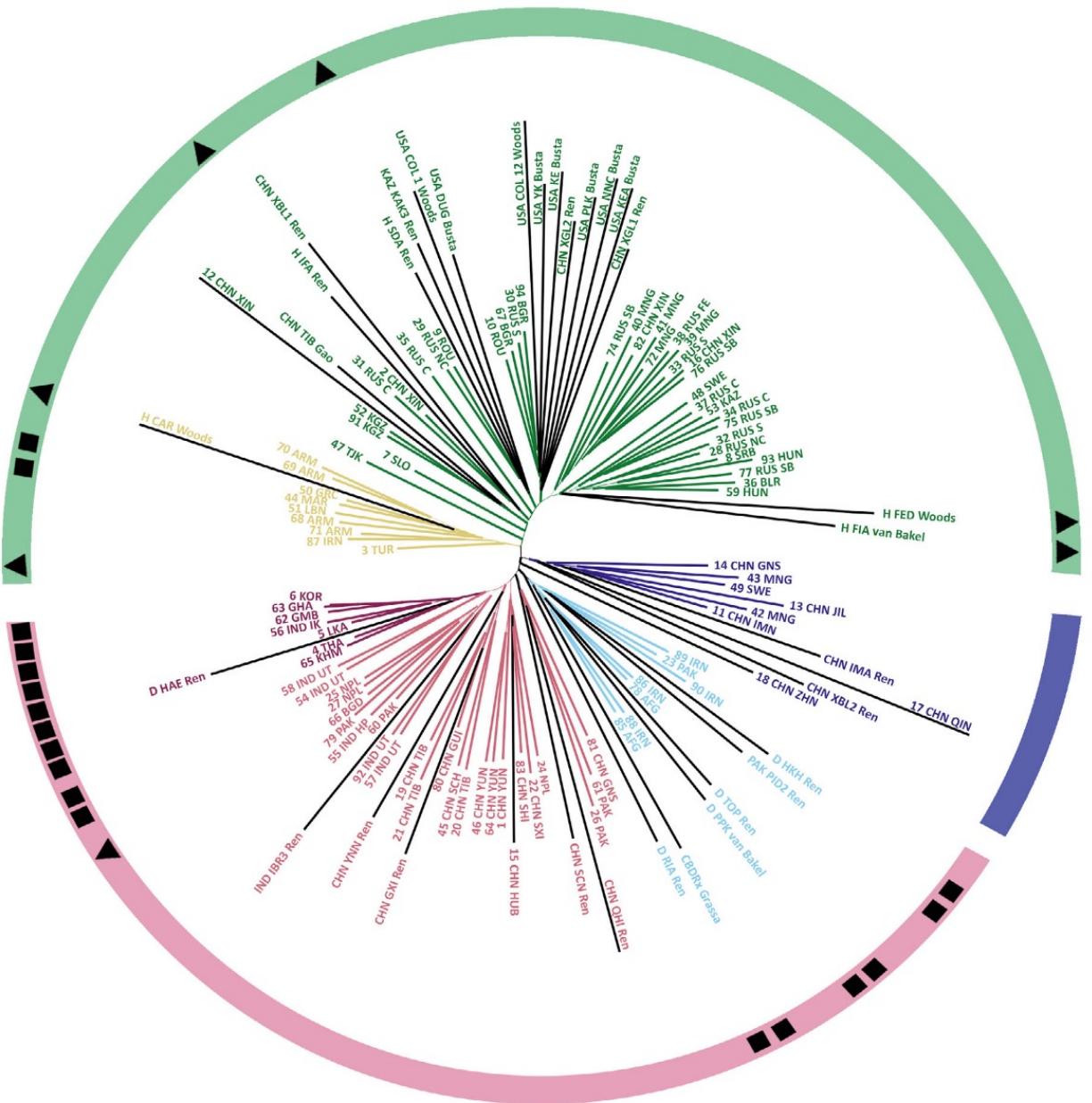


FIGURE 1 Geographic distribution of samples included in this study, with individuals coloured according to the subgroups obtained in the phylogenomic analysis (see Figure 2). The shapes indicate cannabis accession types, them being, wild-growing (circles), fibre/seed (triangles) and drug (squares) types. Additionally, filled shapes are newly analysed Hyb-Seq samples, while empty shapes are NCBI sequence read archive (SRA) corresponding to WGS data mined for our Hyb-Seq targets. The inset shows United States wild-growing populations mined from NCBI SRAs. Drug cultivars mined are not shown. For more detailed information, see Table S1. The map was made with *Natural Earth* (free vector and raster map data @ naturalearthdata.com).





| Phylogenetic subgroups | Phylogenetic groups | Use type and domestication status |
|----------------------------|---------------------|-----------------------------------|
| N China and Mongolia | E Asia | Wild-growing individuals |
| Iranian Plateau | Boreal | ▲ Fiber/seed landrace/cultivars |
| C & S China and Himalayas | Paleotropis | ■ Drug landrace |
| Indoafrica | | |
| Caucasus and Mediterranean | | |
| Euroasberia and Mongolia | | |

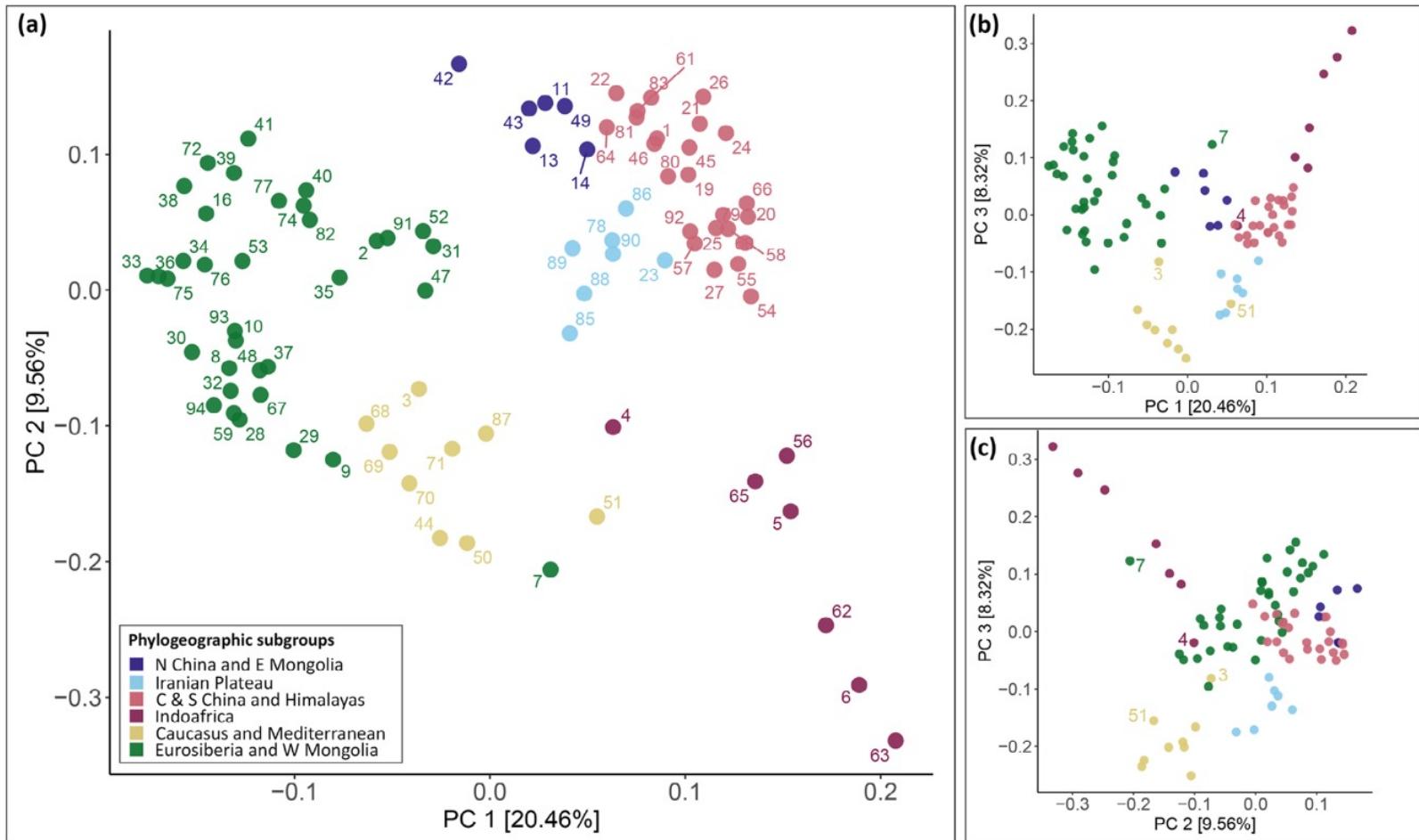


FIGURE 4 Principal component analysis (PCA; performed with PLINK) of *Cannabis* individuals for the 2875 (filtered and unlinked) single nucleotide polymorphisms (SNPs) called from the same 345 nuclear ortholog targets (comprising exons and their flanking regions) used to estimate the nuclear species tree (supercontig data matrix). Colours correspond to the six phylogeographic subgroups identified in the phylogenomic analysis (see Figure 2). (a) First and second PCA axes. (b) First and third PCA axes. (c) Second and third PCA axes.

TABLE 1 Pairwise fixation index (Hudson F_{ST}) values between phylogeographic subgroups.

| Phylogeographic subgroup pairs | | Hudson F_{ST} |
|--------------------------------|-----------------------------|-----------------|
| Eurosiberia and W Mongolia | Indoafrika | 0.155 |
| Caucasus and Mediterranean | Indoafrika | 0.136 |
| N China and E Mongolia | Indoafrika | 0.126 |
| Indoafrika | Iranian plateau | 0.120 |
| Indoafrika | C and S China and Himalayas | 0.090 |
| Caucasus and Mediterranean | N China and E Mongolia | 0.086 |
| Eurosiberia and W Mongolia | C and S China and Himalayas | 0.080 |
| Eurosiberia and W Mongolia | Iranian plateau | 0.077 |
| Caucasus and Mediterranean | C and S China and Himalayas | 0.076 |
| Caucasus and Mediterranean | Eurosiberia and W Mongolia | 0.061 |
| N China and E Mongolia | Iranian plateau | 0.060 |
| N China and E Mongolia | Eurosiberia and W Mongolia | 0.058 |
| Caucasus and Mediterranean | Iranian plateau | 0.048 |
| Iranian plateau | C and S China and Himalayas | 0.039 |
| N China and E Mongolia | C and S China and Himalayas | 0.036 |

Student questions

- They sent their DNA extraction samples to be sequenced by a private company. What are the risks and benefits of doing your own sequence work vs paying it for a company to do it for you?
- How important was human dispersal in the history?
- Potential issues with herbarium specimens?
- What's a paralog? SNP? Exon? Flanking regions? Linkage disequilibrium? F_{ST} ?
- 68,212 SNPs -> 2,875 after filtering. Good idea? Bad?
- In addition, the STRUCTURE model assumes that markers are neutral, which--if they are located within genes--one can reasonably expect that they generally are not. The paper gave some phenomenal background information, but I'm not particularly comfortable with how they analyzed and interpreted their data.
- How might degraded herbarium samples affect results?
- RAXML-NG vs RAXML
- How would you analyze what changed with domestication?
- Why no time calibration?