
1 Ancient DNAs and the Neolithic Chinese super-grandfather Y haplotypes

2
3 Ye Zhang¹, Xiaoyun Lei¹, Hongyao Chen¹, Hui Zhou^{2*} and Shi Huang^{1*}

4
5
6 ¹Center for Medical Genetics, Xiangya Medical School, Central South University,
7 Changsha, Hunan 410078, People's Republic of China

8
9 ²Ancient DNA Laboratory, Research Center for Chinese Frontier Archaeology, Jinlin
10 University, Changchun 130012

11
12 *Corresponding authors:

13 Shi Huang (huangshi@sklmg.edu.cn)

14 Hui Zhou (zhouhui@jlu.edu.cn)

15
16 **Keywords:** Yangshao Culture, Miaodigou Culture, Banpo Culture, Hongshan Culture,
17 Xiajiadian Culture, O2a2b1a1a-F5, O2a1c1a1a1a-F11, ancient DNA

18
19 **Running title:** Ancient DNAs and the Neolithic Chinese super-grandfather Y
20 haplotypes

21 22 23 **Abstract**

24 Previous studies identified 3 Neolithic Han Chinese super-grandfather Y
25 haplotypes, O2a2b1a1a-F5, O2a2b1a2a1-F46, and O2a1b1a1a1a-F11, but their
26 relationships with the archaeological and written records remain unexplored. We here
27 report genome wide DNA data for 12 ancient samples (0.02x-1.28x) from China
28 ranging from 6500 to 2500 years before present (YBP). They belonged to 4 different
29 genetic groups, designated as Dashanqian (DSQ) of Xiajiadian Culture in the
30 Northeast, Banpo (BP) of middle Yangshao Culture in the Central West, Zhengzhou
31 Xishan (ZX) of Miaodigou Culture in the Central Plains, and Others. Present day F5
32 samples were closer in autosomal distances to the ZX and DSQ groups while F11, C,
33 O1, and O2 samples were closer to the BP group. We also sequenced the Y
34 chromosome of one of these ancient samples K12 from DSQ and found both K12 and
35 a previously reported ~4000 year old sample MG48 from Northwest China to have the
36 O2a2b1a1a1a2a-F2137 haplotype, belonging to the most prolific branch
37 O2a2b1a1a1-F438 immediately under F5. We further found close relationships
38 between ZX and DSQ and between ZX and ancient M117 Tibetans or present day
39 Southwest Dai Chinese carrying the F5 subtype O2a2b1a1a6, implicating radiations
40 of F5 subtypes from the putative place of F5 origin in ZX. These results are
41 remarkably consistent with archaeological and written records.

42
43

44 Introduction

45 There are numerous data for human activity in China from the time of the
46 Neolithic period to the beginning of written records. There were the Gaomiao and
47 Pengtoushan Culture of ~5800 BC in the South (Hunan), the Jiahu Culture and
48 Peiligang Culture of 7000-5000 BC in the Central Plains (Henan), the Xinglongwa
49 Culture of 6200-5400 BC and later the Hongshan and Xiajiadian Cultures in the
50 Northeast (Inner Mongolia-Liaoning border), the Dadiwan Culture of 5800-5400 BC in
51 Gansu and Western Shaanxi ¹. At 5000 BC to 3000 BC, the Yangshao Culture was
52 the most popular and existed extensively along the Yellow River in China and
53 flourished mainly in the provinces of Henan, Shaanxi and Shanxi with the early period
54 of the Culture mostly found in Shaanxi and the late period in Henan ¹. Elements of the
55 middle to late Yangshao Culture, the Miaodigou Culture, have been found widely in
56 China, including the Hongshan Culture in the Northeast (3500 BC), indicating the
57 broad cultural migration and influence of this Culture ².

58 By analyzing the Y chromosome haplotype patterns, three Neolithic
59 super-grandfather haplotypes have been discovered that together account for ~40%
60 of present day Han Chinese males ³⁻⁵. The expansion dates are estimated 5400 YBP
61 (years before present) for O3a2c1a-M117-F5 (O2a2b1a1a-F5 or F8, ISOGG 2017),
62 6500 YBP for O3a2c1-F46 (O2a2b1a2a1-F46), or 6800 YBP for O3a1c-002611-F11
63 (O2a1b1a1a1a-F11), and these three haplotypes represent 16%, 11%, and 14% of
64 present day Han Chinese males, respectively. Several historical writings on ancient
65 Chinese mention the great leaders/ancestors around the time of 5000 years ago or
66 earlier, including *Fu Xi*, *Yan Emperor* (Yandi), *Huang Emperor* (Huangdi), and *Chi*
67 *You*.

68 It remains unclear how the Neolithic Cultures were related to the
69 super-grandfather haplotypes and how the different Neolithic Cultures were
70 interconnected. We here addressed these questions by analyzing ancient DNA
71 samples from 10 different sites in Central and Northern China. We found evidence of
72 F5 associated autosomes in the Miaodigou Culture in Henan, F11 associated
73 autosomes in the early Yangshao Culture in Banpo, and both F5 haplotype and F5
74 associated autosomes in the Xiajiadian Culture in Inner Mongolia. The results provide
75 a coherent account of information from the relevant fields.

76

77 Results:

78 Relationships among aDNAs in Central and Northern China

79 We collected 12 skeletal remains from 10 sites in Central and Northern China that
80 were 2500-6500 years old (Table 1). DNAs were extracted and sequenced to different
81 degrees of coverage (0.018x-1.27x). SNPs were called using the Hg19 reference
82 genome. We obtained two sets of SNPs, one slow and the other fast, with the slow set
83 informative for phylogenetics and the fast set representing natural selection and
84 maximum saturation as detailed before ⁶⁻¹². We then merged the SNPs of each of the
85 aDNA samples with those of the 1000 genomes project (1kGP) ¹³, and obtained
86 pairwise genetic distances, i.e., mismatch in homozygous (hom) sites for slow SNPs
87 or mismatch in all sites for fast SNPs, between each aDNA sample and all individuals

88 in the 1kGP.

89

90 **Table 1. Information on ancient samples for which we report the nuclear**
91 **sequence data in this study.**

Sample name	Ages (BP)	# SNPs slow	Coverage	Culture	Sites
FQ17	2500	1200	0.113	East Zhou	Fushan Qiaobei, Shanxi
K2	3000	8300	1.2796	Xiajiadian Upper	Dashanqian, Chifeng, Inner Mongolia
K12	3000	2000	0.2071	Xiajiadian Upper	Dashanqian, Chifeng, Inner Mongolia
S1	3920	350	0.043	Xiaoheyuan	Halahaigou, Chifeng, Inner Mongolia
SG2046	3500	1500	0.1668	Xiajiadian Lower	Sanguan, Hebei
DZ14	5000	1980	0.0712	Late Yangshao	Duzhong, Mianchi Henan
BP38	6500	966	0.0556	Yangshao	Banpo, Xian
ZX167	5300	3722	0.0674	Miaodigou	Zhengzhou Xishan, Henan
ZX107	5300	161	0.0179	Miaodigou	Zhengzhou Xishan, Henan
PAB3002	3000	2300	0.2361	Gaotaishan	Pinganbao, Liaoning
ZK22	4200	500	0.0485	Zhukaigou	Zhukaigou, Ordos, Inner Mongolia
H69	3500	1692	0.0304	Early Shang	Xuecun, Zhengzhou, Henan

92

93 Because different ancient samples were sequenced at different coverages, it is
94 unrealistic to use SNP mismatches to infer relationships as there would be few shared
95 SNPs among different samples. As an alternative, we calculated the correlation
96 coefficient R of two genomes in their distance to the East Asian (ASN) samples in
97 1kGP, assuming that different random sampling of a fraction of the whole set of ~15K
98 slow SNPs are roughly equivalent in representing the whole set. Verification of this R
99 correlation approach has been described previously⁶.

100 We tested the correlation of ancient samples within themselves relative to their
101 correlation with present day Han Chinese in order to better infer the relationships
102 among the ancient samples. We obtained for each aDNA R values with each other
103 and with the 211 Han Chinese samples in 1kGP. We then ranked these R values as
104 shown in Table 2, which then allowed us to infer genetic relationships among these
105 aDNAs. We grouped these 12 aDNAs into 4 groups based on being from the same
106 site and being closely related to each other as indicated by correlation rankings. In
107 Table 2, rank values means the rank of a sample on the column among all 223
108 samples in values of correlation to a sample listed on the row, e.g., K2 from the
109 column was ranked 7th among all samples in correlation with K12 from the row. The 4
110 groups were designated as the Dashanqian (DSQ) group including the 2 samples
111 from the Xiajiadian Culture in Dashanqian in Chifeng (K2 and K12) known to be
112 related to the Hongshan Culture¹⁴, one sample S1 from the Xiaoheyuan Culture at the
113 Hala Haigou site in Chifeng related to the Hongshan Culture¹⁵ and one sample FQ17
114 from the East Zhou dynasty at Fushan Qiaobei Shanxi¹⁶; the Banpo (BP) group
115 including a sample BP38 from the middle Yangshao Culture site at Banpo Xian¹⁷, a
116 late Yangshao sample DZ14 from Duzhong Henan¹⁸, one sample SG2046 from
117 Sanguan Hebei known to be influenced by both Yangshao and Xiajiadian¹⁹, one
118 sample ZK22 from Zhukaigou Inner Mongolia²⁰, and one sample H69 from the early

119 Shang site at Xuecun Henan ²¹; the Zhengzhou Xishan (ZX) group including 2
120 samples from the same Zhengzhou Xishan site of Miaodigou Culture in the Central
121 Plains (ZX167 and ZX107) ²²; and finally the Other group including PAB3002 of
122 Gaotaishan Liaoning ²³.

123

124 **Table 2. Ranks in correlation between each pair of ancient samples among 211**
125 **Han Chinese in 1KG and 12 ancient samples.** Rank values listed means the rank of
126 a sample on the column among all 223 samples in values of correlation to a sample
127 listed on the row.

128

	FQ17	K2	K12	S1	SG2046	DZ14	BP38	ZK22	H69	ZX167	ZX107	PAB3002
FQ17		170	173	169	220	126	182	152	78	74	207	221
K2	20		7	90	217	206	209	61	169	69	167	193
K12	184	103		5	222	180	198	174	109	190	192	143
S1	218	215	149		218	188	206	204	106	211	193	196
SG2046	221	219	222	151		41	99	164	166	173	179	213
DZ14	210	222	220	221	10		1	21	5	215	219	215
BP38	208	221	203	182	171	16		187	189	158	210	126
ZK22	169	187	189	180	199	43	191		16	221	202	113
H69	102	217	155	7	214	35	205	18		222	120	191
ZX167	23	172	164	89	161	128	72	208	196		97	190
ZX107	220	213	198	135	204	182	204	197	94	160		205
PAB3002	222	214	151	61	215	136	90	123	133	204	178	

129

130

131 These groupings mostly showed expected relationships based on geophysical
132 distances, for example, the close genetic relationship of the 3 Chifeng samples in the
133 DSQ group. Also, FQ17 in the DSQ group was located in the Y shaped region known
134 to be a common migration route linking the Central Plains with the Northeast ²⁴. A
135 notable exception was ZX167 who had K2 and FQ17 from DSQ group among the
136 most related aDNA samples rather than ZX107 from the same site (K2 and FQ17
137 ranked higher than all other aDNAs in R correlation to ZX167, being at 69th and 74th,
138 respectively). However, ZX107 did have ZX167 as the most related aDNA (the ranking
139 of 97th was the highest among all aDNAs in R correlation to ZX107), consistent with
140 the geographical location.

141 ZX167 was the 172th ranked in relation to K2 and the closest to K2 among non
142 DSQ samples. K2 was the 69th ranked in relation to ZX167 and the closest to ZX167
143 among all aDNAs samples here. This indicated gene flow between DSQ and ZX,
144 which is consistent with the known archaeological findings of a migration of Miaodigou
145 Culture to the Niuheliang site in Chifeng Inner Mongolia ².

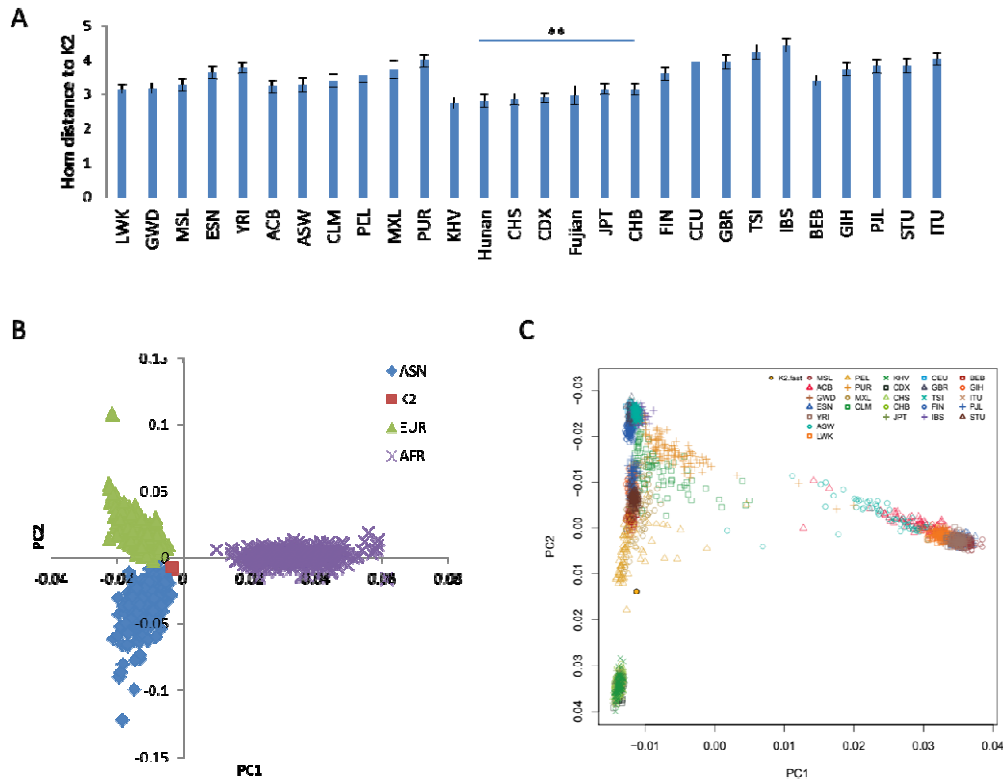
146 DZ14 was the first ranked among all in relation to BP38 with R values
147 substantially higher than the next (0.38 for DZ14 and 0.24 for the second ranked) and
148 BP38 was 16th ranked among all in R to DZ14. This indicated gene flow between
149 Banpo Shannxi (BP38) of middle Yangshao Culture and Duzhong Henan (DZ14) of

150 late Yangshao Culture, consistent with archaeological findings.

151

152 **Relationships between aDNAs and present day samples carrying the**
153 **super-grandfather haplotypes**

154 We first determined to which present day populations the ancient Chinese
155 samples may be closely related. We used the informative slow SNPs to calculate
156 pairwise genetic distances between each ancient DNA and each of the 1kGP samples
157 as described previously^{6,7}. For samples such as K2 with relatively large number of
158 slow SNPs covered, we were able to perform informative distance analyses and
159 principle components analysis (PCA). The results showed K2 to be closest to East
160 Asians in genetic distances (Figure 1A). K2 was closer to Europeans than most East
161 Asians as shown by PCA plots (Figure 1B), which is consistent with the location of K2
162 at Dashanqian Chifeng Inner Mongolia where East and West admixture are known to
163 have occurred²⁵. Interestingly, K2 was more related to South East Asians such as
164 KHV and Hunan people and least related to Northern Chinese CHB (Han Chinese in
165 Beijing), consistent with the known migration of ancient Southwest people from
166 Gaomiao Culture in Hunan to the Jiahu Culture in Henan and in turn to the Hongshan
167 Culture in the Northeast as indicated by archaeological records (8 pointed star)²⁶. In
168 contrast, PCA plots using fast SNPs clustered all aDNA samples here with PEL
169 (Peruvian in Lima, Peru) and as outliers to the ASN samples (Figure 1C and data not
170 shown), which was reminiscent of the surprising finding with European aDNAs that
171 ancient samples (>2000 years old) were often not the direct ancestors of present day
172 Europeans living in the same area²⁷⁻³⁰. These results from the fast SNPs were in
173 direct conflict with all other lines of data and hence most likely wrong, which in turn
174 provided additional support for the theoretical justification for using slow SNPs in
175 demographic inferences⁶⁻¹².
176



177

178

179 **Figure 1. Relationship of ancient Chinese genomes with present day**

180 **populations.** A. Pairwise genetic distance between K2 sample and 1kGP samples.

181 Distances were mismatches in homozygous sites. Standard error of the mean (SEM)

182 is shown. B. PCA plots (PC1 and PC2) of K2 merged with 1kGP samples using slow

183 SNPs. Only ASN, EUR, and AFR samples of 1kGP are selectively shown. C. PCA

184 plots (PC1 and PC2) of K2 merged with 1kGP samples using fast SNPs.

185

186 We performed correlation analyses using distances to European samples of
 187 1kGP to determine whether the aDNAs here were more related to East Asians (CHS
 188 Southern Han Chinese) or Europeans (CEU, Utah Residents with Northern and
 189 Western European Ancestry). All informative aDNAs with R values greater than
 190 expected by random simulations (>0.02) showed greater affinity to CHS than to CEU
 191 (Table 3). Some samples had too few SNPs to be informative. Therefore, although
 192 the numbers of SNPs here were limited, which would weaken the strength of
 193 correlations, most samples were still informative enough to be able to properly identify
 194 the correct group affiliations of these aDNAs.

195

196 **Table 3. Correlation of ancient samples from China with Han Chinese (CHS) or**
 197 **Europeans (CEU) of 1kGP.** Shown are Spearman correlation R values in distance to
 198 EUR samples in 1kGP. Two samples from Han Chinese of 1kGP were shown as
 199 positive controls. Comparisons were deemed non-informative when R values were
 200 $0.1 - 0$ or $-0.1 - 0$, which happens to be found for randomly scrambled distance

201 values as shown for CHS_HG00478.

202

Sample names	R in EUR distance		
	CHS	CEU	T test
FQ17	0.010	0.039	ni
K2	0.423	0.229	<0.01
K12	0.179	0.165	
S1	0.125	0.080	<0.01
SG2046	0.165	0.152	
DZ14	0.136	0.126	
BP38	0.197	0.126	<0.01
ZK22	0.182	0.073	<0.01
H69	0.158	0.138	<0.05
ZX167	0.306	0.211	<0.01
ZX107	-0.024	0.010	ni
PAB3002	0.114	0.088	<0.01
CHS_HG00478	0.623	0.376	<0.01
CHB_NA18611	0.512	0.380	<0.01
CHS_HG00478 random	-0.008	-0.020	ni

203

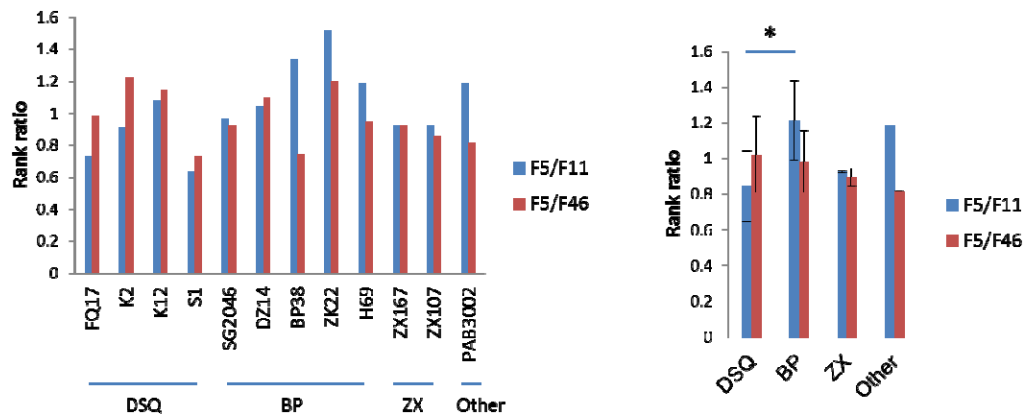
204

205

206 Our previous studies showed that individuals carrying the same Y haplotype were
207 also more related in autosomes, which could only be shown with slow but not fast
208 autosomal SNPs⁷. We next determined the autosomal relationship of each aDNA to
209 present day people carrying the super-grandfather Y haplotypes. We calculated the
210 average rank in R values of each aDNA to the F5 samples relative to the F11 and F46
211 samples in the 211 Han Chinese in 1kGP. A ratio <1 in the F5 rank versus the F11 or
212 F46 rank means closer relationship to the F5 samples relative to the F11 or F46
213 samples. Relative to F11 samples, DSQ and ZX but not BP were found closer to F5
214 samples, indicating that the DSQ and ZX groups had more F5 associated autosomes
215 and less F11 associated autosomes whereas the BP group had the opposite (Figure
216 2A). Relative to the F46 samples, DSQ and BP groups were more evenly related to F5
217 and F46 but the ZX group was still more related to F5, indicating that the ZX group
218 had probably the most F5 associated autosomes among the aDNAs examined here.
219 BP38 was the only sample in the BP group that had significantly more F11 than F46
220 ($P < 0.01$, chi-squared test), indicating close relationship of BP38 with the F11
221 haplotype known to be common in Southern China. While DSQ and ZX groups were
222 both related to F5, DSQ was more related to F46 than ZX.

223 We also did the same analysis using total pairwise distance values calculated
224 from the set of fast evolving SNPs, as would normally be done by the field. The results
225 showed that all 3 aDNA groups (DSQ, DZ, and ZX) had F5/F11 ratio <1 (0.81-0.84),
226 indicating no sensible correlation of any kind. Such dramatic difference in results
227 obtained from the fast and the slow SNPs was consistent with previous findings^{6,7},

228 confirming again that only slow SNPs could be used for informative analyses in
229 phylogenetic studies or in detecting specific association between autosomes and Y
230 haplotypes.



231

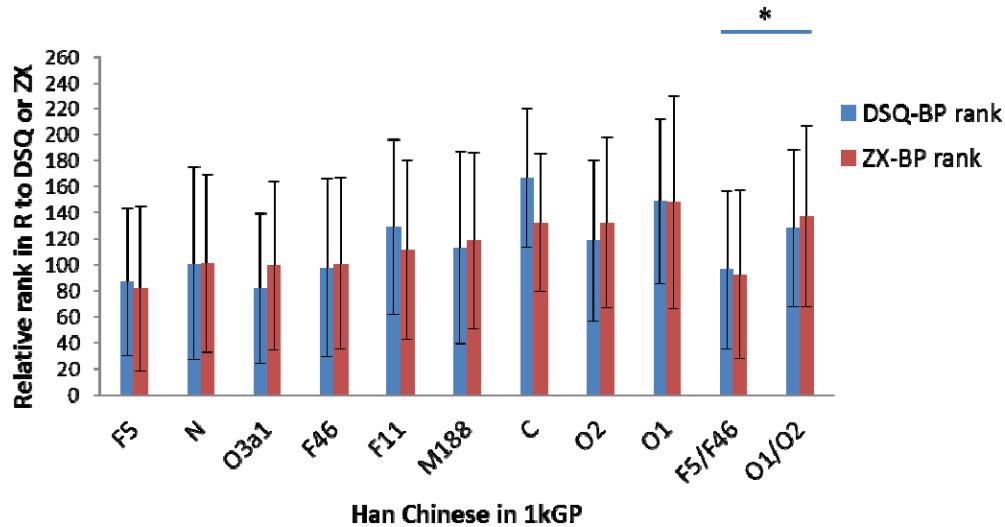
232 **Figure 2. Autosomal relationships of aDNAs with Han Chinese samples of 1kGP**
233 **carrying F5, F11, or F46 haplotypes.** R correlation values with 211 Han Chinese and
234 12 aDNAs were obtained and ranked for each aDNA sample. The average rank to F5,
235 F11, or F46 was calculated for each sample listed here and the ratios in the ranks are
236 shown either for each sample (A) or the average of each group (B).

237

238

239 To determine further the possible Y haplotypes associated with the aDNA groups,
240 we studied the autosomal relationship of aDNAs with present day samples grouped
241 by different haplotypes. We determined the relative distance to DSQ vs BP or to ZX vs
242 BP for different sets of 1kGP samples with each set associated with a particular
243 haplotype. For each sample, we calculated a rank in the subtraction value of ‘the rank
244 to DSQ subtracting the rank to BP (DSQ-BP)’ with small values (<120 with 120 being
245 the middle rank corresponding to the subtraction value 0) meaning higher rank in
246 relation to DSQ relative to BP. We also similarly calculated a rank in the value of ‘the
247 rank to ZX subtracting the rank to BP (ZX-BP)’. The samples of the M134 clade
248 containing both F5 and F46 samples were significantly closer to ZX or DSQ than
249 samples carrying O1 and O2 haplotypes (F5/F46 and O1/O2 in Figure 3). Relative to
250 DSQ, samples carrying F11 and C were the closest to the BP group although not
251 significant. The results suggest that BP group had more autosomal ancestry from
252 Southern China as O1, O2, F11, and C are known to be common in the South.

253



254

255

256

257

258

259

260

261

262

263

264

265

266

267

268

269

270

271

272

273

274

275

276

277

278

279

280

281

282

283

Figure 3. Relative autosomal distance of aDNAs to Han Chinese in 1kGP

grouped by Y haplotypes. Ranks in R to aDNA groups were determined and ranks were shown for the subtraction values of 'rank to DSQ subtracting the rank to BP (DSQ-BP)' or 'rank to ZX subtracting the rank to BP (ZX-BP)'. Small subtraction ranks (<120) means closer rank to DSQ or ZX, and higher value means closer rank to BP.

Ancient DNA relationships with Southwest Chinese carrying M117 or the F5 subtype O2a2b1a1a6

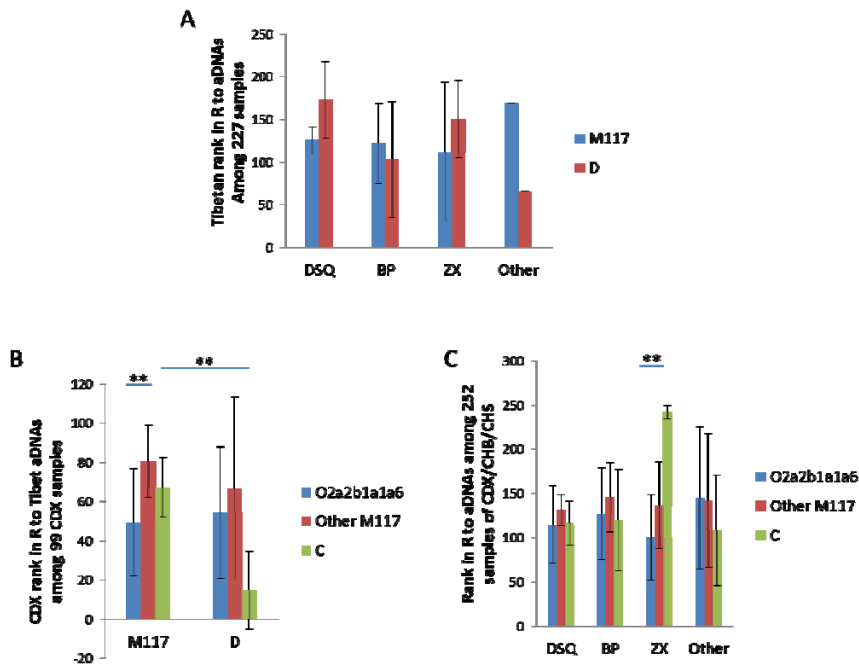
Tibetans are known to commonly carry M117³¹ and a subtype of M117 or F5, O2a2b1a1a6, is common in Southwest Chinese³². If both DSQ and ZX groups may be candidate group for the origin of F5, it would be important to ask which group is more related to Southwest people carrying M117 or O2a2b1a1a6. We made use of 4 previously published Tibetan genomes of 3150-1250 years old, including 3 M117 and one D haplotype (markers downstream of M117 were not covered)³¹. We obtained the slow SNPs set and calculated the pairwise genetic distance with each member of 1kGP for each of the ancient Tibetans.

To determine which of the 12 aDNA samples was closely related to the M117 Tibetans, we obtained correlation R ranking of the 3 M117 Tibetans in relation to each aDNA among 227 samples including 211 Han Chinese, 12 Han aDNAs, and 4 Tibetan aDNAs (Figure 4A). The highest rank for the M117 Tibetans was found in correlation with the ZX group among the 4 aDNA groups. The results indicated possible gene flow between the ZX group and the M117 Tibetans.

The haplotype O2a2b1a1a6 was commonly found in Southwest Chinese such as the Dai³². To further confirm the putative role of ZX in the origin of this haplotype, we studied the CDX samples (Chinese Dai in Xishuangbanna) of the 1kGP. We first confirmed that CDX samples carrying O2a2b1a1a6 were closely related to M117 Tibetans in autosomal genetic distances. We obtained R values of each aDNA with the CDX samples using their genetic distance to East Asian samples in 1kGP. We classified the M117 samples in CDX into 2 different haplogroups (O2a2b1a1a6 and

284 other M117), and calculated the average R of the CDX in relation to the Tibetans. The
285 results showed that CDX samples carrying O2a2b1a1a6 were significantly closer to
286 M117 Tibetans than CDX samples carrying other M117 haplotypes. CDX samples
287 carrying the C haplotype were significantly closer to the D type Tibetan than to M117
288 Tibetans, consistent with C and D belonging to the ABCDE clade⁷. The results
289 indicated that the autosomes of ancient M117 Tibetans were more related to those
290 associated with O2a2b1a1a6, suggesting that ancient Tibetans likely carried the
291 O2a2b1a1a6 haplotype.

292 We next divided the CDX samples into 3 Y groups (O2a2b1a1a6, Other M117,
293 and C) and tested their relationships with the 4 groups of aDNAs here. We obtained
294 correlation values of each aDNA to each of 99 CDX, 103 CHB, and 50 CHS samples
295 of 1kGP. Among the 4 aDNA groups, the ZX group was found closest to the
296 O2a2b1a1a6 samples and the only aDNA group that showed significantly closer
297 distance to O2a2b1a1a6 samples than to C samples (Figure 4C). The results suggest
298 a specific association of the ZX group with populations that carried the O2a2b1a1a6
299 haplotype, indicating either that ZX was enriched with the O2a2b1a1a6 haplotype or
300 that ZX was the ancestral population giving rise to the group carrying the
301 O2a2b1a1a6 haplotype. Given the present day distribution of O2a2b1a1a6 in the
302 Southwest, it is more likely for ZX, being located in the Central Plains, to be the
303 ancestral group to the O2a2b1a1a6 haplotype.



304

305

306 **Figure 4. Relationship of the O2a2b1a1a6 haplotype with aDNA samples. A.**

307 Relationship between ancient Tibetans and the aDNA groups in this study. B.

308 Relationship of Tibetans with CDX samples carrying different Y haplotypes. C.

309 Relationship with aDNA groups of this study for CDX samples carrying different Y

310 haplotypes.

311

312 **Analyses of ancient Y chromosomes**

313 Given the autosomal relationships of both the DSQ and ZX group with present
314 day people carrying the F5 haplotype, it may be expected that at least some members
315 of these two groups may carry the F5 haplotype. We performed Y chromosome
316 sequence analysis on the only sample here that was good enough for high coverage
317 sequencing analysis, the ~3000 year old K12 from Upper Xiajiabian Culture in West
318 Liao River Valley, Northeast China¹⁴. We also analyzed the published sequence of
319 the MG48 individual from the ~4000 year old Mogou site of the Qijia Culture in
320 Northwest China³³. The results showed both K12 and MG48 to be F438-F2137 with
321 no mutations in sites that define downstream haplotypes under F2137. The results
322 provided additional support for the presence of the F5 lineage in the DSQ group.

323

324 **Discussion:**

325 Our results here showed that aDNAs in Central and Northern China from this
326 study could be separated into 4 groups based on autosomal relationships among
327 them. Two among these, ZX and DSQ, were related to F5 associated autosomes. The
328 BP group was more related to autosomes associated with the O1, O2, C, and F11 that
329 are commonly found in the South while the ZX group was more associated with the F5
330 and F46 haplotypes common in the Central Plains and the North (8/12 F11 samples in
331 Han Chinese were CHS and 6/8 F11 in CHS were from Hunan). This indicates that the
332 BP group might be migrants from the South. Consistently, analyses of human skulls of
333 the Duzhong and Banpo sites indicated close relationship with populations from South
334 China^{34,35}. Human skulls from the Zhengzhou Xishan site or other Miaodigou sites
335 such as Shanxian Henan however showed mixed features related to both Yangshao
336 and Dawenkou people^{36,37}. Thus, people from different Miaodigou sites in Henan,
337 such as Duzhong and Xishan, appear to have different cranial features, consistent
338 with DNA findings here. The DNA results confirm the suggestion based on
339 archaeological and historical records that the early Yangshao Culture and its possible
340 predecessor the Peiligang/Jiahu Culture may be associated with migrants under the
341 legendary *Yan* or *Fuxi* Emperor from the South such as the Pengtoushan and
342 Gaomiao Culture^{26,38}.

343 The DSQ group had one sample K12 carrying F5 while the ZX group was
344 non-informative for Y chromosome. There are at least 7 branches immediately under
345 F5³². The F438 branch appears to be of high social economic status (SES) based on
346 it having shorter branch length, more descendant branches, and higher fitness (lower
347 risk for autism)⁴. The O2a2b1a1a1a2a2 haplotype of K12 and MG48 belongs to F438.
348 If the original F5 haplotype had a fitness advantage that might have contributed to its
349 super-grandfather status in the first place, haplotypes with fewer random variations
350 from the original F5 haplotype should be expected to retain the most of the fitness
351 advantage of F5 and as such to be more likely to confer high SES status and produce
352 more descendants. Thus populations in ancient times near the time of the original F5
353 would be expected to be enriched with haplotypes closest to F5. Thus the finding of
354 two of two informative samples that have high coverage sequence data being of the

355 F438-F2137 haplotype is consistent with *a priori* expectation of higher prevalence of a
356 high SES haplotype in ancient times.

357 It appears that the ZX group may be more directly linked to the origin of F5 as it
358 was both related to DSQ in the Northeast and the Southwest people carrying the F5
359 subtype O2a2b1a1a6. The common presence of F438-F2137 in the 3000-4000 YBP
360 time period in the North indicates unlikely the presence of F* haplotype in the North in
361 ancient times. The most parsimonious explanation for these observations is the
362 diversification and radiation of F5 sub-branches from a centrally located population
363 such as ZX where F* might originate.

364 Present day Chinese are thought to be the descendants of *Yan* and *Huang*.
365 Based on archaeological and historical records, scholars have suspected an
366 association of *Yan* with the Yangshao Culture in the Central Plains (but with ancestry
367 from the South such as the Gaomiao Culture in Hunan) and *Huang* with the Hongshan
368 Culture in the Northeast^{24,38,39}.

369 The Miaodigou Culture was a most popular Culture of its time and known to have
370 impacted the Northeast HongShan Culture (and the subsequent Xiajiadian Culture to
371 which the DSQ group belonged), more so than any other Culture of the time such as
372 the Dawenkou Culture². Our DNA findings here suggest that there were people in the
373 Central Plains closely related to the super-ancestor F5 lineage and possibly
374 associated with the Miaodigou Culture. These conclusions from DNA studies are
375 consistent with the suggestion from archaeological and historical studies that the
376 Miaodigou Culture, and in particular the first walled town (made of rammed earth) of
377 Xishan, may be linked to the lineage of *Huang* who is known to be the first to have
378 built walled towns in Chinese history²². People of this lineage are believed to have
379 also lived in the Northeast (Hongshan and Xiajiadian Culture) including the great
380 Zhuangxi Emperor, a grandson of *Huang*, and to have migrated down to the Central
381 Plains in later times during a cold climate period^{24,39}. Samples from the Niuhejiang
382 site of Hongshan Culture are 13.7% for haplotype O2a2b1-M134 (downstream sites
383 remain to be determined) and future studies of more samples are needed to
384 determine if F5 was present at this site⁴⁰.

385 Overall, our study identified the presence of F5 genomes in ancient samples from
386 the Central Plains and the Northeast and implicates the origin of the F5 lineage in the
387 Central Plains and subsequent diversification and migration to the Northeast and
388 Southwest. The remarkable unification of ancient DNA results with archaeological and
389 written records can only be found when we used slow SNPs but not fast SNPs. This
390 provides further confirmation of our new molecular methodology in demographic
391 inferences.

392

393 **Acknowledgments:**

394 This work was supported by the National Natural Science Foundation of China grant
395 81171880 (SH) and 31371266 (H.Z.) the National Basic Research Program of China
396 grant 2011CB51001 (S.H.).

397

398 **Author contributions:**

399 Y.Z., X.L., and H.C. performed experiments and data analysis. H.Z and S.H.
400 conceived and supervised the study. YZ and SH and analyzed the data and wrote the
401 first draft of the manuscript and all authors participated in revising the manuscript.

402

403 **Competing Interests**

404 The authors declare that they have no competing interests.

405

406 **Materials and Methods:**

407

408 **Ancient DNA sequencing**

409 Archaeological sites and samples were described in details in supplementary
410 materials. Two teeth from each sample were collected for DNA sequencing analyses
411 as we did in a previous study³³. Appropriate precautions were taken to ensure the
412 authenticity of the ancient DNA results.

413

414 **Sequence download**

415 We downloaded ancient and modern human genome sequences from the relevant
416 websites using publically available accession numbers.

417

418 **Selection and identification of SNPs**

419 Selecting SNPs: The random selection of autosomal 255K SNPs as fast evolving
420 SNPs (none from the X chr) and the selection of the slow evolving SNPs were as
421 described previously⁷.

422 Calling SNPs: We used publically available software SAMTOOLS, GATK, and
423 VCFTOOLS to call SNPs from either downloaded BAM files or BAM files we
424 generated based on downloaded fastq data or our own data by using
425 Burrows-Wheeler Aligner or BWA software⁴¹⁻⁴³. For better accuracy of calling SNPs,
426 the Phred quality score (Q score) was set at 50, DP \geq 10, and alt frequency \geq 0.8.
427 For the ancient Y chromosome DNAs, the filters for calling SNPs included DP $>$ 4 and
428 alt frequency \geq 0.8.

429

430 **Genetic distance analyses and other population genetics methods**

431 Genetic distance: Pairwise genetic distances of either hom distances (mismatch in
432 hom sites) or total distances (mismatch in het and hom sites combined) were
433 calculated using the custom made software “dist” as previously described⁷. This
434 software is freely available at <https://github.com/health1987/dist> and has been
435 described in detail in previous publications^{44,45}.

436

437 **PC analysis:** We utilized GCTA to analyze data in the PLINK binary PED format to
438 generate two files (*.eigenvec and *.eigenva). We drew PCA plot using *.eigenvec file
439 ^{46,47}.

440

441 **Other methods:** Other common statistical methods used were Student’s t test, chi
442 square test, and Fisher’s exact test, 2 tailed, and Spearman correlation coefficient R

443 analysis using Prism6 of Graphpad.

444

445 **Accession codes**

446 The raw reads of ancient DNAs reported here have been deposited at NCBI with
447 accession number xxxx.

448

449

450 **References:**

451

452 1 Chang, K.-C., Xu, P. & Lu, L. *The formation of Chinese Civilization, an archaeological*
453 *perspective*. (Yale University Press, 2005).

454 2 Guo, D. Why the Hongshan culture is "taproot system" of Chinese ancient culture? *J. Liaoning*
455 *Normal University (Social Science Edition)* **39**, 121-131 (2016).

456 3 Yan, S. *et al.* Y chromosomes of 40% Chinese descend from three Neolithic
457 super-grandfathers. *PLoS ONE* **9**, e105691, doi:10.1371/journal.pone.0105691 (2014).

458 4 He, P. *et al.* Neolithic super-grandfather Y haplotypes, their related surnames, and autism
459 spectrum disorder. *bioRxiv* <https://doi.org/10.1101/077222> (2017).

460 5 Wen, S.-Q., Tong, X.-Z. & Li, H. Y-chromosome-based genetic pattern in East Asia affected by
461 Neolithic transition. *Quaternary International* (2016).

462 6 Yuan, D. & Huang, S. On the peopling of the Americas: molecular evidence for the
463 Paleoamerican and the Solutrean models. *bioRxiv* **bioRxiv 130989**; doi:
464 <https://doi.org/10.1101/130989> (2017).

465 7 Yuan, D. *et al.* Modern human origins: multiregional evolution of autosomes and East Asia
466 origin of Y and mtDNA. *bioRxiv*, doi: <https://doi.org/10.1101/106864> (2017).

467 8 Huang, S. Primate phylogeny: molecular evidence for a pongid clade excluding humans and a
468 prosimian clade containing tarsiers. *Sci China Life Sci* **55**, 709-725 (2012).

469 9 Huang, S. New thoughts on an old riddle: What determines genetic diversity within and
470 between species? *Genomics* **108**, 3-10, doi:10.1016/j.ygeno.2016.01.008 (2016).

471 10 Hu, T. *et al.* The genetic equidistance result, misreading by the molecular clock and neutral
472 theory and reinterpretation nearly half of a century later. *Sci China Life Sci* **56**, 254-261
473 (2013).

474 11 Huang, S. Inverse relationship between genetic diversity and epigenetic complexity. *Preprint*
475 *available at Nature Precedings* doi.org/10.1038/npre.2009.1751.1032 (2009).

476 12 Lei, X., Yuan, J., Zhu, Z. & Huang, S. Collective effects of common SNPs and risk prediction in
477 lung cancer. *Heredity*, doi:10.1038/s41437-41018-40063-41434 (2018).

478 13 Auton, A. *et al.* A global reference for human genetic variation. *Nature* **526**, 68-74,
479 doi:10.1038/nature15393 (2015).

480 14 Peng, S., Zhu, Y., Guo, Z. & Wang, L. Excavation in 1998 on the Dashanqian site in Harqin
481 Banner, Inner Mongolia. *Archaeology* **3**, 31-39 (2004).

482 15 Zhang, Y. Preliminary report on the excavation to the Neolithic cemetery at Hala Haigou
483 village, Chifeng city, Inner Mongolia. *Archaeology* **2**, 19-35 (2010).

484 16 Tian, J. *et al.* The graves of Shang and Zhou dynasties in Qiaobei, Fushan. *J. Ancient*
485 *Civilization* **5**, 347-394 (2006).

486 17 Qian, Y. Study on the burial ways of early Yangshao Banpo site. *Archaeology of West China* **8**,

-
- 487 56-81 (2014).
- 488 18 Wu, Z. Excavation report of 2006 on the Duzhong site, Mianchi, Henan. *Huaxia Archaeology* **3**,
489 3-18 (2006).
- 490 19 Chen, P. A review of lower Xiajiadian culture. *Beijing Cultural Relics and Archaeology* **3**, 55-64
491 (2002).
- 492 20 Mongolia, T. The Zhukaigou site in Inner Mongolia. *Acta Archaeologica Sinica* **3**, 301-332
493 (1988).
- 494 21 Sun, L., Chu, X. & Zhu, H. Study of the human skulls from the early Shang site at Xuecun,
495 Xingyang, Henan. *Huaxia Archaeology* **1**, 55-64 (2013).
- 496 22 Yang, Z. On the character of the city site of the late Yangshao culture at Xishan, Zhengzhou.
497 *Huaxia Archaeology* **1**, 55-59 (1997).
- 498 23 Zhu, H. The report on human bones identification of Pingganbao site. *Acta Archaeologica*
499 *Sinica* **4**, 473-479 (1992).
- 500 24 Guo, D. From the crossroads to Y-type culture zone. *Archaeology of Inner Mongolia* **2**, 96-104
501 (2006).
- 502 25 Liu, S. *et al.* Genomic Analyses from Non-invasive Prenatal Testing Reveal Genetic
503 Associations, Patterns of Viral Infections, and Chinese Population History. *Cell* **175**, 347-359
504 e314, doi:10.1016/j.cell.2018.08.016 (2018).
- 505 26 He, G. *Prehistoric West Hunan and Classical Chinese Historical Legends*. (2013).
- 506 27 Fu, Q. *et al.* The genetic history of Ice Age Europe. *Nature* **534**, 200-205,
507 doi:10.1038/nature17993 (2016).
- 508 28 Lazaridis, I. *et al.* Ancient human genomes suggest three ancestral populations for
509 present-day Europeans. *Nature* **513**, 409-413, doi:10.1038/nature13673 (2014).
- 510 29 Mallick, S. *et al.* The Simons Genome Diversity Project: 300 genomes from 142 diverse
511 populations. *Nature* **538**, 201-206, doi:10.1038/nature18964 (2016).
- 512 30 Mathieson, I. *et al.* Genome-wide patterns of selection in 230 ancient Eurasians. *Nature* **528**,
513 499-503, doi:10.1038/nature16152 (2015).
- 514 31 Jeong, C. *et al.* Long-term genetic stability and a high-altitude East Asian origin for the
515 peoples of the high valleys of the Himalayan arc. *Proc Natl Acad Sci U S A* **113**, 7485-7490,
516 doi:10.1073/pnas.1520844113 (2016).
- 517 32 Poznik, G. D. *et al.* Punctuated bursts in human male demography inferred from 1,244
518 worldwide Y-chromosome sequences. *Nat Genet* **48**, 593-599, doi:10.1038/ng.3559 (2016).
- 519 33 Li, J. *et al.* Ancient DNA reveals genetic connections between early Di-Qiang and Han Chinese.
520 *BMC Evol Biol* **17**, 239, doi:10.1186/s12862-017-1082-0 (2017).
- 521 34 Sun, L. & Wu, Z. Study on the human skeletal samples of the late Yangshao culture at
522 Duzhong Mianchi. *Huaxia Archaeology* **3**, 100-109 (2010).
- 523 35 Han, K. X. Some issues in the study of the characteristics of anthropology materials of the
524 Neolithic Yangshao culture. *Prehistory Research* **3**, 240-256 (1988).
- 525 36 Han, K. & Q., P. A study of the human skeletal remains unearthed from the tombs of the
526 Miaodigou II culture at Shanxian, Henan province. *Acta Archaeologica Sinica* **2**, 255-270
527 (1979).
- 528 37 Wei, D., Zhang, Y. & Zhu, H. Research on the human skeletal remains from Xishan, Zhengzhou.
529 *Cultural Relics of Central China* **2**, 111-119 (2015).
- 530 38 Xu, S. Z. Restudy of the Huangdi Period--the start of Chinese Civilization. *Archaeology and*

- 531 *Cultural Relics* **4**, 19-26 (1997).
- 532 39 Lei, G. Z. Historical remains, Shanhaijing, and Hongshan Culture: the legend of HuangDi and
533 ZhuanXu and the Hongshan Culture. *Journal of Liaoning Teacher College (Social Science*
534 *Edition)* **47**, 115-118 (2006).
- 535 40 Cui, Y. *et al.* Y Chromosome analysis of prehistoric human populations in the West Liao River
536 Valley, Northeast China. *BMC Evol Biol* **13**, 216, doi:10.1186/1471-2148-13-216 (2013).
- 537 41 Li, H. *et al.* The Sequence Alignment/Map format and SAMtools. *Bioinformatics* **25**,
538 2078-2079, doi:10.1093/bioinformatics/btp352 (2009).
- 539 42 Danecek, P. *et al.* The variant call format and VCFtools. *Bioinformatics* **27**, 2156-2158,
540 doi:10.1093/bioinformatics/btr330 (2011).
- 541 43 McKenna, A. *et al.* The Genome Analysis Toolkit: a MapReduce framework for analyzing
542 next-generation DNA sequencing data. *Genome Res* **20**, 1297-1303,
543 doi:10.1101/gr.107524.110 (2010).
- 544 44 Zhu, Z., Lu, X., Yuan, D. & Huang, S. Close genetic relationships between a spousal pair with
545 autism-affected children and high minor allele content in cases in autism-associated SNPs.
546 *Genomics*, 10.1016/j.ygeno.2016.1012.1001, doi:10.1016/j.ygeno.2016.12.001 (2016).
- 547 45 Zhu, Z. *et al.* Collective effects of SNPs on transgenerational inheritance in *Caenorhabditis*
548 *elegans* and budding yeast. *Genomics* **106**, 23-29, doi:10.1016/j.ygeno.2015.04.002 (2015).
- 549 46 Yang, J., Lee, S. H., Goddard, M. E. & Visscher, P. M. GCTA: a tool for genome-wide complex
550 trait analysis. *Am J Hum Genet* **88**, 76-82, doi:10.1016/j.ajhg.2010.11.011 (2011).
- 551 47 Purcell, S. *et al.* PLINK: A tool set for whole-genome association and population-based linkage
552 analyses. *Am J Hum Genet* **81**, 559-575 (2007).

553

554

555

556 **Tables:**

557 **Table 1. Information on ancient samples for which we report the nuclear**
558 **sequence data in this study.**

Sample name	Ages (BP)	# SNPs slow	Coverage	Culture	Sites
FQ17	2500	1200	0.113	East Zhou	Fushan Qiaobei, Shanxi
K2	3000	8300	1.2796	Xiajiadian Upper	Dashanqian, Chifeng, Inner Mongolia
K12	3000	2000	0.2071	Xiajiadian Upper	Dashanqian, Chifeng, Inner Mongolia
S1	3920	350	0.043	Xiaohayan	Halahaigou, Chifeng, Inner Mongolia
SG2046	3500	1500	0.1668	Xiajiadian Lower	Sanguan, Hebei
DZ14	5000	1980	0.0712	Late Yangshao	Duzhong, Mianchi Henan
BP38	6500	966	0.0556	Yangshao	Banpo, Xian
ZX167	5300	3722	0.0674	Miaodigou	Zhengzhou Xishan, Henan
ZX107	5300	161	0.0179	Miaodigou	Zhengzhou Xishan, Henan
PAB3002	3000	2300	0.2361	Gaotaishan	Pinganbao, Liaoning
ZK22	4200	500	0.0485	Zhukaigou	Zhukaigou, Ordos, Inner Mongolia
H69	3500	1692	0.0304	Early Shang	Xuecun, Zhengzhou, Henan

559

560

561

562 **Table 2. Ranks in correlation between each pair of ancient samples among 211**
 563 **Han Chinese in 1KG and 12 ancient samples.** Rank values means the rank of a
 564 sample on the column among all 223 samples in values of correlation to a sample
 565 listed on the row.
 566

	FQ17	K2	K12	S1	SG2046	DZ14	BP38	ZK22	H69	ZX167	ZX107	PAB3002
FQ17		170	173	169	220	126	182	152	78	74	207	221
K2	20		7	90	217	206	209	61	169	69	167	193
K12	184	103		5	222	180	198	174	109	190	192	143
S1	218	215	149		218	188	206	204	106	211	193	196
SG2046	221	219	222	151		41	99	164	166	173	179	213
DZ14	210	222	220	221	10		1	21	5	215	219	215
BP38	208	221	203	182	171	16		187	189	158	210	126
ZK22	169	187	189	180	199	43	191		16	221	202	113
H69	102	217	155	7	214	35	205	18		222	120	191
ZX167	23	172	164	89	161	128	72	208	196		97	190
ZX107	220	213	198	135	204	182	204	197	94	160		205
PAB3002	222	214	151	61	215	136	90	123	133	204	178	

567

568

569 **Table 3. Correlation of ancient samples from China with Han Chinese (CHS) or**
 570 **Europeans (CEU) of 1kGP.** Shown are Spearman correlation R values in distance to
 571 EUR samples in 1kGP. Two samples from Han Chinese of 1kGP were shown as
 572 positive controls. Comparisons were deemed non-informative when R values were
 573 0.1 – 0 or -0.1 – 0, which happens to be found for randomly scrambled distance
 574 values as shown for CHS_HG00478.
 575

Sample names	R in EUR distance		
	CHS	CEU	T test
FQ17	0.010	0.039	ni
K2	0.423	0.229	<0.01
K12	0.179	0.165	
S1	0.125	0.080	<0.01
SG2046	0.165	0.152	
DZ14	0.136	0.126	
BP38	0.197	0.126	<0.01
ZK22	0.182	0.073	<0.01
H69	0.158	0.138	<0.05
ZX167	0.306	0.211	<0.01
ZX107	-0.024	0.010	ni
PAB3002	0.114	0.088	<0.01
CHS_HG00478	0.623	0.376	<0.01
CHB_NA18611	0.512	0.380	<0.01
CHS_HG00478 random	-0.008	-0.020	ni

576

577

578 **Figure Legends:**

579 **Figure 1. Relationship of ancient Chinese genomes with present day**

580 **populations.** A. Pairwise genetic distance between K2 sample and 1kGP samples.

581 Distances were mismatches in homozygous sites. SEMs are shown. B. PCA plots

582 (PC1 and PC2) of K2 merged with 1kGP samples using slow SNPs. Only ASN, EUR,

583 and AFR samples of 1kGP are selectively shown. C. PCA plots (PC1 and PC2) of K2

584 merged with 1kGP samples using fast SNPs.

585

586 **Figure 2. Autosomal relationships of aDNAs with Han Chinese samples of 1kGP**

587 **carrying F5, F11, or F46 haplotypes.** R correlation values with 211 Han Chinese and

588 12 aDNAs were obtained and ranked for each aDNA sample. The average rank to F5,

589 F11, or F46 was calculated for each sample listed here and the ratios in the ranks are

590 shown either for each sample (A) or the average of each group (B).

591

592 **Figure 3. Relative autosomal distance of aDNAs to Han Chinese in 1kGP**

593 **grouped by Y haplotypes.** Ranks in R to aDNA groups were determined and ranks

594 were shown for the subtraction value of 'rank to DSQ subtracting the rank to BP

595 (DSQ-BP)' or 'rank to ZX subtracting the rank to BP (ZX-BP)'. Small subtraction ranks

596 means closer rank to DSQ or ZX, and higher value means closer rank to BP.

597

598 **Figure 4. Relationship of the O2a2b1a1a6 haplotype with aDNA samples.** A.

599 Relationship between ancient Tibetans and the aDNA groups in this study. B.

600 Relationship of Tibetans with CDX samples carrying different Y haplotypes. C.

601 Relationship with aDNA groups of this study for CDX samples carrying different Y

602 haplotypes.

603