1    **Sexual dimorphism of facial width-to-height ratio in human skulls and faces: A meta-**

2    **analytical approach**

3

4    Robin S. S. Kramer*

5    Department of Psychology, University of York, York, YO10 5DD, UK

6

7    * Corresponding author.

8    *E-mail address:* remarknibor@gmail.com (R. S. S. Kramer).

9

10    Word count: 5891

11

12

20

21

22    **Abstract**

23

24         Facial width-to-height ratio (FWHR), defined as the width of the face divided by the

25    upper facial height, is a cue to behaviour. Explanations for this link often involve the idea

26    that FWHR is sexually dimorphic, resulting from intersexual selection pressures. However,

27    few studies have considered sexual dimorphism in skulls since the original paper on this

28    topic, and it is possible that different explanations may be required if faces show sex

29    differences but skulls do not. Here, meta-analyses of skulls found that men did have larger

30    FWHR than women, although this effect was small. However, after categorising samples by

31    ethnicity and geographical origin, meta-analyses only found evidence of sex differences in

32    East Asians, and again, this effect was small. A re-analysis of previous studies after

33    excluding skull samples found little evidence of sexual dimorphism in faces. Again,

34    considering ethnicities separately, I found no differences for White samples but a medium-

35    sized effect with East Asians, although this was not statistically significant with only three

36    samples. Taken together, I found no reason to consider FWHR as a sexually dimorphic

37    measure in skulls or faces, at least not universally, and so accounts based upon this

38    assumption need rethinking if researchers are to explain the relationship between FWHR and

39    behaviour.

40

41    *Keywords:* Facial width-to-height ratio; Sexual dimorphism; Skulls; Faces; Testosterone

42

43    **1. Introduction**

44

45         The idea that the human face provides social information is not a new one (Darwin,

46    1872). We can determine the identity (Bruce & Young, 1986), sex (Burton et al., 1993), age

47  (Rhodes, 2009), and ethnicity (Montepare & Opeyo, 2002) of a stranger with relative ease, as

48  well as more dynamic and changing information like emotional state (Elfenbein & Ambady,

49  2002). There is also evidence that trait information like personality, physical and mental

50  health, and even sexual orientation can be perceived with some accuracy from faces alone

51  (Jones et al., 2012; Kramer & Ward, 2010; Rule et al., 2009; Scott et al., 2013).

52      In 2007, researchers provided evidence of one particular facial measure, the width-to-

53  height ratio (FWHR – see Fig.1; Weston et al., 2007), which they found to be sexually

54  dimorphic in human skulls, and has since been the subject of intense investigation as a cue to

55  numerous behaviours. While overall size differences play a large role in general skull

56  dimorphism (Calcagno, 1981; Lestrel, 1974; Rightmire, 1970), Weston and colleagues

57  suggested that this ratio difference was instead due to developmental differences in shape

58  trajectories during puberty. Specifically, the height of the upper face (defined as the nasion-

59  prosthion distance) in adults is similar in men and women, while the (bizygomatic) width is

60  larger in men. In other words, while sex differences in skulls are expected simply because

61  men grow to be larger than women, the bizygomatic width in males shows additional growth

62  at puberty beyond this predicted increase. The researchers argued that this difference in skull

63  shape might result from intersexual selection pressures, so that a region of the face has

64  evolved which highlights the distinction between men and women.

65      Why evidence of sexual dimorphism predicts an association between FWHR and

66  behaviour is less clear. If female preferences led to increased facial width in men (although

67  evidence actually suggests that wider faces are judged to be less attractive; Geniole et al.,

68  2015), it may not necessarily follow that within-sex differences are correlated with

69  behaviours. More intuitively, intrasexual selection pressures (e.g., male-male competition)

70  could have resulted in increased success for wider-faced men, resulting in an appearance–

71  behaviour link, especially if these two factors have the same underlying mechanism

72 (testosterone, for instance). Might both explanations overlap, whereby a facial cue that

73 highlights 'maleness' to women has become associated with masculinity (both in appearance

74 and behaviour) in men? Of course, there is no reason to assume that the mechanisms

75 underlying sex differences in facial development are the same as those that may drive a

76 within-sex association between appearance and behaviour.

77 While the precise account and its relationship with sexual dimorphism remains unclear,

78 FWHR does appear to function as a social cue. Levels of masculine characteristics (e.g.,

79 aggression, dominance, deception) in men correlate with FWHR, as do perceived levels of

80 these traits (for meta-analyses, see Geniole et al., 2015; Haselhuhn et al., 2015). The

81 explanation for this FWHR–behaviour association is thought to involve testosterone (Carré &

82 McCormick, 2008; Sell et al., 2009), which may influence both facial development and

83 behavioural characteristics. Indeed, initial research found significant associations between

84 FWHR in men and baseline levels of testosterone, as well as testosterone changes in response

85 to potential mate exposure (Lefevre et al., 2013).

86 Somewhat problematically for this account, FWHR may not actually be sexually

87 dimorphic in faces (Kramer et al., 2012; Lefevre et al., 2012; Özener, 2012) or skulls

88 (Gómez-Valdés et al., 2013; Stirrat et al., 2012). Of course, it may be that different

89 mechanisms drive facial development in men and women, allowing for testosterone-produced

90 correlates of behaviour in men without differences between the sexes (Lefevre et al., 2013).

91 In a recent meta-analysis of this field, the authors found significant (but small) sex

92 differences when considering studies of both skulls and faces together (Geniole et al., 2015),

93 as well as for subsets of studies (2D photographs versus other materials). However, it is not

94 clear whether differences remain when only skulls are analysed since this distinction was not

95 made in their analyses. It may be that skulls do not show sex differences in FWHR but faces

96 do, perhaps through evolved cues that utilise soft tissue deposits, which differ in men and

97  women (Enlow, 1982). This would be an important caveat when investigating the

98  explanatory mechanisms linking behaviour and facial measures.

99      One potential issue with previous investigations is that they have not considered

100  populations separately based upon ethnicity or geographical origin. Given evidence of

101  between-population differences in skulls (Gill & Rhine, 2004; İşcan & Steyn, 1999; Ousley et

102  al., 2009), the inclusion of all groups into a single analysis will inherently suffer from this

103  additional source of noise. It may be that FWHR dimorphism is present in some

104  ethnicities/populations but not others, and this could account for the mixed results that have

105  previously been found with faces. This would also be an important caveat for theories of

106  dimorphism and signalling.

107      The other problem for the 'FWHR–testosterone–behaviour' account is that FWHR may

108  not actually be associated with testosterone. In recent research investigating several samples

109  and reporting a combined meta-analysis, no relationship was found between FWHR in adult

110  men and baseline testosterone or competition-induced testosterone reactivity (Bird et al.,

111  2016). Even during adolescence, when testosterone is hypothesised to impact facial growth

112  (Weston et al., 2007), no relationship was found between male FWHR and testosterone levels

113  or other known testosterone-derived traits (Hodges-Simeon et al., 2016). Indeed, FWHR

114  showed no change during adolescence and no growth spurt, contrary to predictions.

115      In the current work, I focussed specifically on whether FWHR is sexually dimorphic in

116  adult human skulls using a meta-analytical approach. Given that the popular topic of FWHR

117  as an important facial cue originated from this initial finding (Weston et al., 2007), it is worth

118  further examination using multiple large samples. I also considered geographical and ethnic

119  origins as potential factors in order to allow for the likelihood that populations may differ.

120  For this reason, I revisit the topic of FWHR sex differences in faces, again considering

121  ethnicity as a potential influence. Importantly, prior large-scale analyses in this area have yet

122 to consider the distinction between faces and skulls, and the possibility (and evolutionary

123 implications) that there may be FWHR sex differences in one but not the other.

124

125 **2. Methods**

126

127 *2.1. Previous research*

128

129       All peer-reviewed and published manuscripts that investigated human skull FWHR

130 separately for men and women were included. This involved searching through all articles

131 that cited Weston et al. (2007), the first paper to propose this measure as a topic of interest.

132 Conveniently, all articles prior to the end of 2014 had already been identified in the recent

133 meta-analysis by Geniole et al. (2015), and no newer research (as of May 2016) or omissions

134 were found. This resulted in the inclusion of three peer-reviewed manuscripts.

135       In total, these publications described eight separate skull databases: six (Gómez-Valdés

136 et al., 2013), one (Stirrat et al., 2012), and one (Weston et al., 2007). Problematically, the

137 authors reported, and the previous meta-analysis utilised, summary database values for

138 FWHR rather than separating these into specific populations in terms of origin/ethnicity. For

139 example, Stirrat and colleagues provided means and standard deviations for their full sample,

140 which included a mixture of White and non-White skulls. Similarly, Gómez-Valdés and

141 colleagues reported the average FWHR dimorphism for each database, which did not allow

142 for the analysis of separate populations, incorporating different sample sizes, etc. For

143 instance, the Howells (1973, 1989, 1995) database alone contained 26 groups (of varying

144 sizes) originating from almost as many countries.

145       To address this issue, I contacted the authors (Gómez-Valdés et al., 2013) and obtained

146 summary statistics for their databases, separately for each population. This would allow the

147 calculation of an effect size for each group rather than each database. Unfortunately, the

148 authors were unable to provide data regarding two of their previously reported databases

149 (Hallstat and Mexico City Penitentiary) due to ethical constraints and data property issues,

150 and so these two sets were not included in the current meta-analysis. In addition, several

151 populations were removed before analyses because there was substantial overlap across their

152 databases. For example, the Ourga specimens appeared in both the 2D and Pucciarelli sets.

153 Whenever multiple occurrences were found, the repeated case with the smaller sample size

154 was removed. This was because the second appearance often featured fewer specimens and

155 so was assumed to be a subset of the larger sample, and this was confirmed by the authors

156 through correspondence.

157 For Weston's original sample (Weston et al., 2007), the set contained individuals from

158 several different southern African populations. However, each of these was not represented in

159 sufficient numbers to allow separation into subgroups, and the authors reported previous

160 work demonstrating that these populations were comparable (de Villiers, 1968). I therefore

161 considered this set as a single population for the purposes of analysis.

162 Finally, in order to allow for populations as separate sets/studies, I incorporated into the

163 meta-analysis only groups that included a minimum of two men and two women. Additional

164 data/populations were discarded.

165

166 *2.2. Databases*

167

168 Although Gómez-Valdés et al. (2013) provided their summary data regarding the

169 Howells database (Howells, 1973, 1989, 1995), I was able to obtain this set independently

170 (see Section 2.2.1). I therefore used my own calculated values for these populations, given

171 that it was preferable to work with the raw data when available. Similarly for the database

172 used by Stirrat et al. (2012), I obtained the original set independently (see Section 2.2.2).

173 From these data, I was able to calculate summary statistics separately for each population.

174       In addition, I obtained four new skull databases in order to increase the number of

175 populations included in the meta-analysis and improve the reliability of the findings. A full

176 summary of the final databases included in the analysis can be found in the Supplementary

177 Materials.

178

179 *2.2.1. William W. Howells craniometric data set*

180

181       This database contained information on a large number of specimens. All skulls were

182 from adults (approximately 18 years old and above, as determined by dental development,

183 although exact age was not known), and sex and origin were included. The skulls were

184 pooled from historical collections from various institutions internationally, and contained 30

185 indigenous populations. A full description can be found in Howells' monographs (1973,

186 1989, 1995). FWHR was calculated using the bizygomatic breadth and nasion-prosthion

187 height, measured directly from the skulls. The usable data included here comprised 2412

188 individuals from 26 populations.

189

190 *2.2.2. Database for forensic anthropology in the United States, 1962-1991*

191

192       This forensic database was created in order to represent the ethnic diversity and

193 demographic structure of the US population. A full description can be found in the codebook

194 (Jantz & Moore-Jansen, 2006). From the initial set, specimens were excluded due to missing

195 cranial measures, if they were aged below 18, or if their sex or race were not reported. Given

196 that the current analysis relies heavily upon the accuracy of these two pieces of information, I

197 also chose to exclude specimens where there was a label for sex/race but the researchers had

198 specified uncertainty in their reporting of these categories. The usable data included here

199 comprised 665 individuals from two populations.

200

201 *2.2.3. Hamann-Todd human osteological collection*

202

203       This database contained skulls collected in Cleveland, Ohio, and the surrounding area

204 in the first half of the twentieth century, and housed at the Cleveland Museum of Natural

205 History. The sex, age, and ethnicities of cadavers were recorded, along with skull

206 measurements. Those specimens where age, sex, or the necessary measures were missing,

207 were excluded from analyses, along with any specimens younger than 18 years old. As

208 above, FWHR was calculated using the bizygomatic breadth and nasion-prosthion height,

209 measured directly from the skulls. The usable data included here comprised 2614 individuals

210 from three populations.

211

212 *2.2.4. Forensic 3D database*

213

214       A database of specimens was created in order to facilitate forensic identification using

215 geometric morphometrics. The set incorporates skulls taken from several different

216 collections, including the Roger J. Terry anatomical skeletal collection (Smithsonian

217 Institution, Washington, DC) and the forensic data bank (University of Tennessee,

218 Knoxville), and features populations from around the world. All specimens were adults over

219 the age of 18 (determined based on standard growth and development). Sex and origin

220 information was available, along with three-dimensional coordinates for craniofacial

221     landmarks. These were used to calculate bizygomatic breadth and nasion-prosthion height.

222     The usable data included here comprised 419 individuals from six populations.

223

224     *2.2.5. "Hispanic" populations craniometric database*

225

226         This database was created by the North Carolina State University's forensic analysis

227     lab in order to investigate sources of admixture in "Hispanic" populations, and represents

228     individuals from European, South and Central American countries. Sex, age, and information

229     regarding skull width and height, were missing for many of the items. In addition, specimens

230     under the age of 18 (determined based on standard growth and development) were also

231     excluded. Finally, all specimens originating from Peru were removed since a subset of these

232     also appeared in the Howells database. The usable data included here comprised 50

233     individuals from two populations.

234

235     *2.2.6. Modern, cranial, postcranial and dental metrics database*

236

237         The majority of these measurements were recorded by Peter Brown at Australian

238     National University. The sample consists of Australian Aborigines, Southern Chinese, and

239     European skulls. Many of the specimens have missing data, including age (although all were

240     classed as adults). Bizygomatic breadth and nasion-prosthion height were measured directly

241     from the skulls. The usable data included here comprised 259 individuals from five

242     populations.

243

244     *2.3. Ethnicity*

245

246    Populations were broadly categorised where sufficient numbers were present. I make

247    no claims about the nature of race as a useful, or even justifiable, biological concept

248    (Smedley & Smedley, 2005). Instead, simply in order to investigate whether populations may

249    show different patterns based on their similar ethnic or geographical origins, I used broad

250    umbrella terms that may help to highlight commonalities and differences in skull sexual

251    dimorphism. For example, I used 'White' to incorporate individuals originating from Europe

252    and North America, where their origin was considered to be Caucasian. Similarly, 'Black'

253    included both North American and African skulls that had been previously classified as

254    Black, African American or Native African. I also used the broad categories 'East Asian',

255    'Australian Aboriginals', 'Pacific Islands', and 'South American'.

256    It is important to note that craniometric variation at the global level (between

257    geographic regions or populations) is much lower than within local populations (Relethford,

258    2002). Indeed, perhaps as low as 11-14% of global diversity exists between regions, where

259    the rest falls within regions (Relethford, 1994). Therefore, it would be preferable to

260    categorise the current samples using far narrower groupings than the ones presented here,

261    focussing on true populations (e.g., Han Chinese) rather than more general regions (East

262    Asian), as this would likely improve the chances of finding group-level differences.

263    Unfortunately, the availability of samples prevents such narrow categorisations while still

264    maintaining a reasonable subgroup size. However, such within-subgroup noise means that

265    any statistical differences between groups are perhaps even more suggestive of

266    ethnic/population differences.

267

268    *2.4. Statistics*

269

270       All data were based upon differences between men and women, and so I chose to use

271     Cohen's *d* as the effect size. Analyses were carried out using customised Microsoft Excel

272     spreadsheets, based upon suggestions outlined in previous work (Geniole et al., 2015), as

273     well as the formulae and guidelines provided by Cumming (2012). Specifically, the pooled

274     estimate of the standard deviation within groups was used as the standardiser for *d*. In

275     addition, unbiased estimates of $\delta$ ($d_{unb}$) were used in all cases after applying Hedges'

276     adjustment to *d* to account for small samples [$d*(1-(3/((4*df) - 1)))$ where *df* is degrees of

277     freedom]. Finally, the effect size for each dataset was weighted by the inverse of its variance

278     before calculation of the mean weighted effect size.

279       The 95% confidence intervals for each study's effect size (see Supplementary

280     Materials) were calculated (using Wilsons's online calculator:

281     http://www.campbellcollaboration.org/escalc/) around *d* rather than $d_{unb}$ because these

282     provide a better estimate of the intervals around the population value, $\delta$ (Cumming, 2012).

283       All analyses presented here use random effects models, which assume that the

284     population means estimated by the different studies are randomly chosen from a

285     superpopulation (heterogeneity). Fixed effects models, in contrast, assume that every study

286     estimates the same mean (homogeneity), and so any variation in sample effects is due to

287     sampling error alone. Although random effects models are more complex, they are also

288     considered more realistic and are generally recommended (Cumming, 2012). One further

289     advantage is that fixed effects models are simply a special case of random effects models,

290     where the population variance happens to be zero (Hunter & Schmidt, 2004). As discussed

291     below, I also found statistical evidence to suggest that the samples were heterogeneous.

292

293     **3. Results**

294

295     *3.1. Meta-analysis of skulls for all populations*

296

297     The Supplementary Materials provides a summary of the eleven databases included in

298     the meta-analysis, comprising 4918 men and 2924 women from 87 populations.

299

300     *3.1.1. Heterogeneity*

301

302     The two previous meta-analyses in this field disagreed with regard to which type of

303     model was most suitable: fixed (Haselhuhn et al., 2015) or random effects (Geniole et al.,

304     2015). As such, I first discuss the evidence supporting the use of random effects models here.

305     Several statistics were considered in order to explore the heterogeneity of the databases

306     (whether different samples estimate different population effect sizes or a single one). First, I

307     found statistically significant variability between study means, $Q(86) = 162.01$, $p < .0001$. In

308     other words, the observed variation across studies (162) was greater than the expected

309     variation (which is equal to the degrees of freedom, 86). However, this test can be both poor

310     at detecting true heterogeneity due to low power, and can have excessive power with

311     many/larger studies (Higgins et al., 2003). As such, other measures (e.g., $T^2$ or $I^2$) often prove

312     more informative with regard to the amount of inconsistency, but can also allow comparisons

313     to be made across analyses. The estimated variance of the true effect sizes (the amount of true

314     heterogeneity) appears relatively low ($T^2 = 0.05$), although notice that this means our estimate

315     of their *SD* is 0.22 while the mean effect size itself (see Section 3.1.2) is only 0.09 in the

316     same units. Further, about half ($I^2 = 46.92\%$, considered moderately large) of this observed

317     variance is real, i.e., due to heterogeneity rather than simply being spurious. Finally, the

318     'diamond ratio' (Cumming & Calin-Jageman, 2017), calculated by dividing the margin of

319     error produced by the random effects model by the one given by the fixed effects model, was

320 1.59. Since this is a ratio, a value of 1 would suggest little heterogeneity, and so the current

321 value implies there is heterogeneity present. Taken together, there is evidence here to proceed

322 with random effects models, which indeed many recommend the use of in all situations

323 (Cumming, 2012).

324

325 *3.1.2. Results*

326

327       The results of the meta-analysis found that men's FWHR was slightly larger than

328 women's, $N = 7941$, $k = 87$, mean weighted $d = 0.09$, 95% CI [0.01, 0.17], $p = .022$. This

329 result is in line with the previous findings of Geniole et al. (2015), whose effect size was

330 0.11, [0.03, 0.20], and included studies measuring both skulls and faces.

331       Inspection of the 87 samples (see S1 Fig) identified eight apparent outliers, which had

332 effect sizes with confidence intervals that did not overlap with those of the mean weighted

333 effect size. (These are noted in the Supplementary Materials.) Excluding these outlying effect

334 sizes increased the mean weighted effect size and decreased the confidence interval, $N =$

335 5955, $k = 79$, mean weighted $d = 0.10$, [0.04, 0.16], $p = .002$. In addition, the variability

336 between study means was no longer significant, $Q(78) = 96.93$, $p = .072$. Finally, the

337 variance due to heterogeneity could now be considered small, $I^2 = 19.53\%$.

338       As noted above, 47% of the variation across samples was due to heterogeneity rather

339 than chance (prior to the removal of outliers). Given this degree of variability, it may be the

340 case that one or more study-level characteristics (moderators) could account for some of this

341 variation. First, I consider ethnicity as a potential moderator.

342

343 *3.2. Meta-analyses of skulls by ethnicity*

344

345    *3.2.1. Subgroup analysis*

346

347         I carried out a subgroup analysis to investigate whether ethnicity was a moderator in the

348    overall meta-analysis. Similar to a conventional analysis of variance, this method allows for

349    the comparison of effect sizes across subgroups (here, ethnicities) in order to determine the

350    effect of group-level variables. Study samples were labelled as one of six broad categories of

351    ethnicity/origin, excluding the remaining populations that did not fall within one of these

352    categories. A summary of these can be seen in Table 1.

353         For random effects models, I need to estimate the value of $\tau^2$, the variance of true effect

354    sizes across the set of studies/samples. Since I am interested in estimating the mean and

355    sampling distribution for each subgroup, I need an estimate of $\tau^2$ within each subgroup. There

356    is no *a priori* reason to assume that the true study-to-study dispersion is the same within all

357    subgroups, and so I use a separate estimate of $\tau^2$ for each subgroup. However, if there are

358    only a few studies within subgroups (e.g., fewer than five; see Table 1) then the estimates of

359    $\tau^2$ are likely to be imprecise. In such cases, the recommendation is to use a pooled estimate in

360    order to increase accuracy (Borenstein et al., 2009). Here, I present the results of both

361    methods.

362         Using random effects with separate estimates of $\tau^2$, I found that ethnicity was not a

363    statistically significant moderator, $Q_{between}(5) = 7.51$, $p = .186$. Utilising a pooled estimate of

364    $\tau^2$ produced a similar result, $Q_{between}(5) = 4.80$, $p = .440$. However, the ability to demonstrate

365    that ethnicity is a moderator in these analyses requires large variation between subgroup

366    means and little variation within subgroups. Problematically, at least some of the subgroups

367    show large within-group variation (see Table 1), making any moderator effects difficult to

368    detect.

369     One way to address this large within-subgroup variation is to remove any outlying

370     effect sizes. As with the overall meta-analysis (Section 3.1.2), subgroups were inspected, this

371     time comparing effect sizes to the mean weighted effect size for that particular subgroup.

372     Only two samples were excluded, one from the Black subgroup (Weston et al., 2007) and one

373     from the South America subgroup. Subgroup analyses were then repeated. Using random

374     effects with separate estimates of $\tau^2$, I found that ethnicity was not a statistically significant

375     moderator, $Q_{between}(5) = 9.77$, $p = .082$, although the result is approaching significance.

376     Utilising a pooled estimate of $\tau^2$ produced a similar result, $Q_{between}(5) = 7.49$, $p = .187$.

377     Given these results, I carried out a meta-analysis for each subgroup in order to

378     investigate ethnicity further, acknowledging that formal tests were only suggestive of a

379     moderating effect but failed to reach statistical significance.

380

381     *3.2.2. Separate meta-analyses for ethnicities*

382

383     For each of the six broad categories of ethnicity/origin, I carried out a separate meta-

384     analysis. The results are summarised in Table 1.

385

386     **Table 1.** A summary of the meta-analyses for the six categories.

| Category | N | k | Q | $I^2$ (%) | d | 95% CI | p |
|---|---|---|---|---|---|---|---|
| White | 2849 | 9 | 23.42* | 65.84 | -0.07 | [-0.24, 0.11] | .450 |
| Pacific Islands | 404 | 4 | 10.05* | 70.15 | 0.05 | [-0.31, 0.41] | .798 |
| Black | 1951 | 9 | 24.64* | 67.53 | 0.06 | [-0.14, 0.26] | .572 |
| South America | 696 | 25 | 42.17* | 43.08 | 0.13 | [-0.08, 0.35] | .222 |
| Australian Aboriginals | 238 | 3 | 2.63 | 23.90 | 0.13 | [-0.16, 0.42] | .386 |
| East Asia | 487 | 8 | 6.38 | 0 | 0.26 | [0.09, 0.43] | .002 |

387     $N$ is the total sample size, k is the number of studies, or populations in this case, and $d$ is the

388     mean weighted effect size. $Q$ and $I^2$ are measures related to the amount of heterogeneity in

389     the group. * Significant at an alpha level of 0.05. Negative values of $I^2$ are set to zero

390     (Higgins et al., 2003).

391

392         While the meta-analysis of all populations of skulls suggested that men have larger

393     FWHR than women (although effect sizes less than 0.2 are considered small), the separate

394     analyses for each ethnicity and geographic origin perhaps support a different interpretation.

395     Only East Asian skulls show evidence of an effect (and even then, it is small), with no other

396     categories suggesting sexual dimorphism. In fact, if these eight East Asian populations were

397     excluded, the remaining populations as a whole no longer provide (statistically significant)

398     evidence of a sex difference, $N = 7454$, $k = 79$, mean weighted $d = 0.07$, 95% CI [-0.01,

399     0.15], $p = .087$.

400         For completeness, the Black and South America subgroups were re-analysed after

401     exclusion of the previously mentioned outliers (Section 3.2.1.). For the Black populations, the

402     result remained qualitatively unchanged, $N = 1891$, $k = 8$, mean weighted $d = -0.02$, 95% CI

403     [-0.19, 0.15], $p = .818$. For the South American populations, removal of the outlier produced

404     an almost significant result, $N = 680$, $k = 24$, mean weighted $d = 0.18$, [-0.01, 0.37], $p = .067$.

405         Interestingly, many of the categories show moderate to large inconsistencies ($I^2$),

406     perhaps suggesting the presence of further moderators or simply that the broad labels used

407     here remain too inclusive and require additional subdivisions (see Section 2.3). In contrast,

408     the East Asian studies showed no observed heterogeneity (and, as a result, provide values

409     similar to a fixed effects model). This suggests that these particular samples are all measuring

410     the same construct.

411

412    *3.3. Reanalysis of Geniole et al. (2015)*

413

414         In a previous meta-analysis, Geniole et al. (2015) found a small but statistically

415    significant difference between the FWHR of men and women. However, the researchers

416    included studies conducted on both faces and skulls, and did not discuss the possibility of

417    differences between ethnicities. Previous evidence has shown that facial dimensions vary

418    across ethnicities (Fang et al., 2011). I therefore reanalysed their data while taking into

419    account the possibility of differences between faces and skulls, and the potential effect of

420    ethnicity.

421

422    *3.3.1. Skulls versus faces*

423

424         The meta-analysis results for samples of faces, using the authors' unaltered data (Table

425    S1 from Geniole et al., 2015) but *excluding* the eight samples which investigated skulls,

426    found no (statistically significant) evidence of the presence of sex differences, $N = 4161$, $k =$

427    24, mean weighted $d = 0.12$, 95% CI [-0.01, 0.25], $p = .068$. Of course, although no longer

428    significantly different from zero, this result remains similar to the original analysis with skull

429    samples included, mean weighted $d = 0.11$, [0.03, 0.20]. Indeed, a subgroup analysis using

430    random effects with separate estimates of $\tau^2$ showed that source (skulls, faces) was not a

431    statistically significant moderator, $Q_{\text{between}}(2) = 0.17$, $p = .920$. However, this analysis

432    included samples of all ethnicities, which may be one reason why I found a large amount of

433    within-group heterogeneity ($I^2$): skulls – 47%, faces – 75%.

434

435    *3.3.2. Meta-analyses of White faces*

436

437        For several of the studies of faces in the previous meta-analysis, the ethnicities of the

438        participants were either unreported in the original papers or included a mixed sample. A

439        meta-analysis of the studies with only White samples (the only category which included more

440        than two studies) again found no evidence of the presence of sex differences in faces, $N =$

441        2037, $k = 9$, mean weighted $d = $ -0.05, [-0.21, 0.11], $p = .559$. In this case, the result is very

442        similar to that of White skulls (see Table 1).

443

444    *3.4. Meta-analysis of Chinese faces*

445

446        In skulls, only East Asians seemed to provide some evidence of larger FWHR in men

447        (see Section 3.2.2). However, the majority of studies on FWHR in faces have focussed on

448        White populations. Therefore, in order to investigate whether faces of this ethnicity

449        demonstrate sex differences in FWHR, I carried out a meta-analysis that included two

450        previous studies of Korean faces (Huh, 2013; Huh et al., 2014) and an additional sample of

451        images that I had collected a few years ago. Front-facing photographs, with neutral

452        expressions, were taken of 135 Chinese students (56 men; age range 19-44; age $M = 23.15$,

453        $SD = 3.30$) at Bangor University, UK. In line with previous research, images were rotated so

454        that both pupils were aligned to the same transverse plane, and then FWHR was calculated

455        using the horizontal distance between the zygions, and the vertical distance between the

456        highest point of the upper lip and the highest point of the eyelids (Kramer et al., 2012).

457        A summary of the three samples can be found in the Supplementary Materials. A meta-

458        analysis of these face samples found no (statistically significant) evidence of the presence of

459        sex differences, $N = 331$, $k = 3$, mean weighted $d = 0.44$, [-0.24, 1.12], $p = .204$. However,

460        the point estimate (considered around a medium-sized effect) is noticeably higher than for the

461 White faces, and the confidence intervals include large effect sizes as well as zero. Of course,

462 the inclusion of additional samples would provide a more precise estimate of the true effect.

463

464 **4. Discussion**

465

466     The current meta-analyses provide evidence that casts doubt on what seems to be the

467 currently accepted story regarding FWHR dimorphism. Although an overall analysis of

468 human skulls found a very small (though statistically significant) effect, where men showed

469 larger FWHR than women, there is an argument to be made for considering ethnicities

470 separately (Gill & Rhine, 2004; İşcan & Steyn, 1999; Ousley et al., 2009). Subgroup analyses

471 were suggestive of a moderating effect of ethnicity (although these failed to reach statistical

472 significance). However, after carrying out separate analyses for six ethnicities/geographical

473 origins, I found that only East Asians demonstrated FWHR sexual dimorphism, and again,

474 this effect was small, although notably larger than for the other groups. Given the between-

475 population differences in skulls, it may be that some populations show sexual dimorphism for

476 this ratio while others do not. Such a result goes against the idea that FWHR represents an

477 evolved signal as a result of sex differences during development if we accept that differences

478 are only (weakly) present in a limited number of populations.

479     The investigation of ethnicity presented here is necessarily limited by the availability of

480 samples that can be reasonably pooled within subgroups. Given that most craniometric

481 variation exists within local populations rather than between geographic regions (Relethford,

482 1994), it is important that further research investigates differences between specific

483 populations (e.g., Han Chinese) where sufficient numbers of samples can be obtained. In this

484 way, we should be better able to identify which populations demonstrate sexual dimorphism

485 in FWHR and which do not.

486      Of course, it may be that skulls do not show sexual dimorphism with regard to FWHR,

487      and instead, humans have evolved a signalling system based upon facial soft tissue deposits.

488      Sex differences in soft tissue thickness may play a role in FWHR cues (Enlow, 1982; Lefevre

489      et al., 2013), and previous evidence has established an association between FWHR and body

490      mass index (Coetzee et al., 2010; Kramer et al., 2012). A recent meta-analysis found a "small

491      but significant difference" in FWHR across 32 samples (Geniole et al., 2015, p. 14) where the

492      majority of studies involved face measurements. Crucially, the researchers did not carry out a

493      separate analysis for faces alone after excluding samples of skulls. (Also remember that their

494      face and skull samples did not differentiate between ethnicities.) Here, I found that a

495      replication of their analysis after excluding skull samples failed to identify a (statistically

496      significant) difference between men and women. In addition, by controlling for ethnicity and

497      considering only White populations (the ethnicity best represented), the effect decreased

498      further to the point where there was no suggestion of a sex difference. However, mirroring

499      the results with skulls, there may be some evidence suggesting an effect for East Asian faces,

500      although more samples are needed before we can make any firm claims regarding the

501      presence of FWHR dimorphism in specific ethnicities.

502      The first studies in this field provided strong evidence that men had larger FWHR than

503      women for both skulls ($d_{unb} = 0.84$, 95% CI [0.32, 1.38]; Weston et al., 2007) and faces ($d_{unb}$

504      $= 0.50$, [0.10, 0.96]; Carré & McCormick, 2008). In the intervening years, researchers have

505      found mixed results, and I show here that there is no compelling evidence to support the

506      initial hypothesis that FWHR is sexually dimorphic. Interestingly, the sample reported in

507      Weston et al. was identified here as an outlying effect size in both the meta-analysis of all

508      populations and in the Black subgroup analysis. While there is no general evidence of sexual

509      dimorphism in the current work, there may be an exception for specific ethnicities or

510      populations, although even in these cases, the effects remain relatively small. Statistically, it

511 makes little sense to state with complete certainty that there is "no effect" (i.e., no difference

512 between the FWHR of men and women), but I argue that consideration of the evidence

513 presented here leads us to conclude that, at most, the effect is very small or absent.

514 How do we reconcile this conclusion with the growing evidence that FWHR is a

515 reliable predictor of various behaviours (Geniole et al., 2015; Haselhuhn et al., 2015)? If we

516 rule out the idea that FWHR cues are the result of sexual selection, through the exaggeration

517 of a sexually dimorphic trait, then it is still possible that other mechanisms are responsible for

518 the FWHR–behaviour association. However, the most likely contender was a testosterone-

519 based mechanism but this has failed to find recent support in large samples (Bird et al., 2016;

520 Hodges-Simeon et al., 2016).

521 Could we explain facial cues using a perception-based mechanism instead? There is

522 very strong evidence that those with higher FWHR are perceived to be more aggressive and

523 dominant (Geniole et al., 2015). Perhaps this is because relatively wider faces subtly

524 resemble angry expressions, and people's perceptions are the result of an overgeneralisation

525 of their judgements of emotional expressions (Said et al., 2009). Interestingly, angry faces do

526 not actually have higher FWHR than neutral expressions (Kramer, 2016; although Marsh et

527 al., 2014, find the opposite result when FWHR is measured differently).

528 A second possibility has been couched in terms of "babyfaceness" – having a rounder

529 face, and as a result, a higher FWHR. Previously, evidence had shown that boys who

530 appeared more babyfaced displayed higher academic achievement if they were motivated to

531 do so, but if they came from lower socioeconomic backgrounds, they showed more criminal

532 behaviours (Zebrowitz et al., 1998a). In addition, early babyfaceness was associated with

533 assertiveness and hostility later in life (Zebrowitz et al., 1998b). This result was explained as

534 a self-defeating prophecy, whereby boys compensated for the warm or naïve stereotypes that

535 people applied to them by manifesting personality traits that counteracted expectations.

536    However, there is now evidence that childhood babyfaceness and infant FWHR are both

537    associated with infant temperament (Arcus & Kagan, 1995; Zebrowitz et al., 2015). These

538    researchers also found a significant correlation between babyfaceness and adult FWHR.

539    Taken together, the suggestion is that a bolder temperament from a larger FWHR in infancy

540    extends through life, resulting in babyfaced adults (who have larger FWHR) demonstrating

541    more assertive and aggressive behaviours. In support of this idea, longitudinal studies have

542    shown that infant temperament predicts behaviour in adolescence and adulthood (e.g.,

543    Schwartz et al., 2012). While the mechanism linking infant temperament and facial

544    appearance remains unknown, possible candidates include prenatal glucocorticoid or

545    testosterone exposure (Arcus & Kagan, 1995).

546        In conclusion, I find a lack of evidence suggesting FWHR differences between men and

547    women, both in skulls and in faces. Considered alongside recent evidence that FWHR does

548    not appear to be associated with testosterone, researchers should now seek new mechanisms

549    in order to explain the relationship between FWHR and behaviour.

550

555

556    **References**

557

558    Arcus, D., & Kagan, J. (1995). Temperament and craniofacial variation in the first two years.

559        *Child Development, 66*, 1529-1540.

560    Bird, B. M., Cid Jofré, V. S., Geniole, S. N., Welker, K. M., Zilioli, S., Maestripieri, D., …

561         Carré, J. M. (2016). *Evolution and Human Behavior*. Advance online publication.

562    Borenstein, M., Hedges, L. V., Higgins, J. P. T., & Rothstein, H. R. (2009). *Introduction to*

563         *meta-analysis*. Chichester, UK: Wiley.

564    Bruce, V., & Young, A. (1986). Understanding face recognition. *British Journal of*

565         *Psychology, 77*, 305-327.

566    Burton, A. M., Bruce, V., & Dench, N. (1993). What's the difference between men and

567         women? Evidence from facial measurement. *Perception, 22*, 153-176.

568    Calcagno, J. M. (1981). On the applicability of sexing human skeletal material by

569         discriminant function analysis. *Journal of Human Evolution, 10*, 189-198.

570    Carré, J. M., & McCormick, C. M. (2008). In your face: Facial metrics predict aggressive

571         behaviour in the laboratory and in varsity and professional hockey players. *Proceedings*

572         *of the Royal Society B: Biological Sciences, 275*, 2651-2656.

573    Coetzee, V., Chen, J., Perrett, D. I., & Stephen, I. D. (2010). Deciphering faces: Quantifiable

574         visual cues to weight. *Perception, 39*, 51-61.

575    Cumming, G. (2012). *Understanding the new statistics: Effect sizes, confidence intervals,*

576         *and meta-analysis*. New York: Routledge.

577    Cumming, G., & Calin-Jageman, R. (2017). *Introduction to the new statistics: Estimation,*

578         *open science, and beyond*. New York: Routledge.

579    Darwin, C. R. (1872). *The expression of the emotions in man and animals*. London: John

580         Murray.

581    de Villiers, H. (1968). *The skull of the South African Negro: A biometrical and*

582         *morphological study*. Johannesburg: Witwatersrand University Press.

583    Elfenbein, H. A., & Ambady, N. (2002). On the universality and cultural specificity of

584         emotion recognition: A meta-analysis. *Psychological Bulletin, 128*, 203-235.

585      Enlow, D. H. (1982). *Handbook of facial growth*. Philadelphia: Saunders.

586      Fang, F., Clapham, P. J., & Chung, K. C. (2011). A systematic review of inter-ethnic

587          variability in facial dimensions. *Plastic and Reconstructive Surgery, 127*, 874-881.

588      Geniole, S. N., Denson, T. F., Dixson, B. J., Carré, J. M., & McCormick, C. M. (2015).

589          Evidence from meta-analyses of the facial width-to-height ratio as an evolved cue of

590          threat. *PLoS ONE, 10*(7), e0132726.

591      Gill, G. W., & Rhine, S. (Eds.). (2004). *Skeletal attribution of race: Methods for forensic*

592          *anthropology*. Albuquerque, NM: Maxwell Museum of Anthropology.

593      Gómez-Valdés, J., Hünemeier, T., Quinto-Sánchez, M., Paschetta, C., de Azevedo, S.,

594          González, M. F., … González-José, R. (2013). Lack of support for the association

595          between facial shape and aggression: A reappraisal based on a worldwide population

596          genetics perspective. *PLoS ONE, 8*(1), e52317.

597      Haselhuhn, M. P., Ormiston, M. E., & Wong, E. M. (2015). Men's facial width-to-height

598          ratio predicts aggression: A meta-analysis. *PLoS ONE, 10*(4), e0122637.

599      Higgins, J. P. T., Thompson, S. G., Deeks, J. J., & Altman, D. G. (2003). Measuring

600          inconsistency in meta-analyses. *BMJ, 327*, 557-560.

601      Hodges-Simeon, C. R., Hanson Sobraske, K. N., Samore, T., Gurven, M., & Gaulin, S. J. C.

602          (2016). Facial width-to-height ratio (fWHR) is not associated with adolescent

603          testosterone levels. *PLoS ONE, 11*(4), e0153083.

604      Howells, W. W. (1973). *Cranial variation in man: A study by multivariate analysis of*

605          *patterns of differences among recent human populations: Vol. 67. Papers of the*

606          *Peabody Museum of Archaeology and Ethnology*. Cambridge, MA: Harvard University

607          Press.

608    Howells, W. W. (1989). *Skull shapes and the map: Craniometric analyses in the dispersion*

609        *of modern Homo: Vol. 79. Papers of the Peabody Museum of Archaeology and*

610        *Ethnology.* Cambridge, MA: Harvard University Press.

611    Howells, W. W. (1995). *Who's who in skulls: Ethnic identification of crania from*

612        *measurements: Vol. 82. Papers of the Peabody Museum of Archaeology and Ethnology.*

613        Cambridge, MA: Harvard University Press.

614    Huh, H. (2013). Digit ratios, but not facial width-to-height ratios, are associated with the

615        priority placed on attending to faces versus bodies. *Personality and Individual*

616        *Differences, 54*, 133-136.

617    Huh, H., Yi, D., & Zhu, H. (2014). Facial width-to-height ratio and celebrity endorsements.

618        *Personality and Individual Differences, 68*, 43-47.

619    Hunter, J. E., & Schmidt, F. L. (2004). *Methods of meta-analysis: Correcting error and bias*

620        *in research findings* (2nd ed.). Newbury Park, CA: Sage Publications.

621    İşcan, M. Y., & Steyn, M. (1999). Craniometric determination of population affinity in South

622        Africans. *International Journal of Legal Medicine, 112*, 91-97.

623    Jantz, R. J., & Moore-Jansen, P. H. (2006). Database for forensic anthropology in the United

624        States, 1962-1991: Inter-university Consortium for Political Social Research (ICPSR)

625        [distributor]; 2006.

626    Jones, A. L., Kramer, R. S. S., & Ward, R. (2012). Signals of personality and health: The

627        contributions of facial shape, skin texture, and viewing angle. *Journal of Experimental*

628        *Psychology: Human Perception and Performance, 38*, 1353-1361.

629    Kramer, R. S. S. (2016). Within-person variability in men's facial width-to-height ratio.

630        *PeerJ, 4*, e1801.

631    Kramer, R. S. S., Jones, A. L., & Ward, R. (2012). A lack of sexual dimorphism in width-to-

632        height ratio in White European faces using 2D photographs, 3D scans, and

633        anthropometry. *PLoS ONE, 7*(8), e42705.

634    Kramer, R. S. S., & Ward, R. (2010). Internal facial features are signals of personality and

635        health. *The Quarterly Journal of Experimental Psychology, 63*, 2273-2287.

636    Lefevre, C. E., Lewis, G. J., Bates, T. C., Dzhelyova, M., Coetzee, V., Deary, I. J., & Perrett,

637        D. I. (2012). No evidence for sexual dimorphism of facial width-to-height ratio in four

638        large adult samples. *Evolution and Human Behavior, 33*, 623-627.

639    Lefevre, C. E., Lewis, G. J., Perrett, D. I., & Penke, L. (2013). Telling facial metrics: Facial

640        width is associated with testosterone levels in men. *Evolution and Human Behavior, 34*,

641        273-279.

642    Lestrel, P. E. (1974). Some problems in the assessment of morphological size and shape

643        differences. *Yearbook of Physical Anthropology, 18*, 140-162.

644    Marsh, A. A., Cardinale, E. M., Chentsova-Dutton, Y. E., Grossman, M. R., & Krumpos, K.

645        A. (2014). Power play: Expressive mimicry of valid agonistic cues. *Social

646        Psychological and Personality Science, 5*, 684-690.

647    Montepare, J. M., & Opeyo, A. (2002). The relative salience of physiognomic cues in

648        differentiating faces: A methodological tool. *Journal of Nonverbal Behavior, 26*, 43-59.

649    Ousley, S., Jantz, R., & Freid, D. (2009). Understanding race and human variation: Why

650        forensic anthropologists are good at identifying face. *American Journal of Physical

651        Anthropology, 139*, 68-76.

652    Özener, B. (2012). Facial width-to-height ratio in a Turkish population is not sexually

653        dimorphic and is unrelated to aggressive behavior. *Evolution and Human Behavior, 33*,

654        169-173.

655    Relethford, J. H. (1994). Craniometric variation among modern human populations.

656         *American Journal of Physical Anthropology, 95*, 53-62.

657    Relethford, J. H. (2002). Apportionment of global human genetic diversity based on

658         craniometrics and skin color. *American Journal of Physical Anthropology, 118*, 393-

659         398.

660    Rhodes, M. G. (2009). Age estimation of faces: A review. *Applied Cognitive Psychology, 23*,

661         1-12.

662    Rightmire, G. P. (1970). Bushman, Hottentot and South African Negro crania studied by

663         distance and discrimination. *American Journal of Physical Anthropology, 33*, 169-195.

664    Rule, N. O., Ambady, N., & Hallett, K. C. (2009). Female sexual orientation is perceived

665         accurately, rapidly, and automatically from the face and its features. *Journal of*

666         *Experimental Social Psychology, 45*, 1245-1251.

667    Said, C. P., Sebe, N., & Todorov, A. (2009). Structural resemblance to emotional expressions

668         predicts evaluation of emotionally neutral faces. *Emotion, 9*, 260-264.

669    Schwartz, C. E., Kunwar, P. S., Greve, D. N., Kagan, J., Snidman, N. C., & Bloch, R. B.

670         (2012). A phenotype of early infancy predicts reactivity of the amygdala in male adults.

671         *Molecular Psychiatry, 17*, 1042-1050.

672    Scott, N. J., Kramer, R. S. S., Jones, A. L., & Ward, R. (2013). Facial cues to depressive

673         symptoms and their associated personality attributions. *Psychiatry Research, 208*, 47-

674         53.

675    Sell, A., Cosmides, L., Tooby, J., Sznycer, D., von Rueden, C., & Gurven, M. (2009). Human

676         adaptations for the visual assessment of strength and fighting ability from the body and

677         face. *Proceedings of the Royal Society B: Biological Sciences, 276*, 575-584.

678    Smedley, A., & Smedley, B. D. (2005). Race as biology is fiction, racism as a social problem

679        is real: Anthropological and historical perspectives on the social construction of race.

680        *American Psychologist, 60*, 16-26.

681    Stirrat, M., Stulp, G., & Pollet, T. V. (2012). Male facial width is associated with death by

682        contact violence: Narrow-faced males are more likely to die from contact violence.

683        *Evolution and Human Behavior, 33*, 551-556.

684    Weston, E. M., Friday, A. E., & Liò, P. (2007). Biometric evidence that sexual selection has

685        shaped the hominin face. *PLoS ONE, 2*(8), e710.

686    Zebrowitz, L. A., Andreoletti, C., Collins, M. A., Lee, S. Y., & Blumenthal, J. (1998a).

687        Bright, bad, babyfaced boys: Appearance stereotypes do not always yield self-fulfilling

688        prophecy effects. *Journal of Personality and Social Psychology, 75*, 1300-1320.

689    Zebrowitz, L. A., Collins, M. A., & Dutta, R. (1998b). The relationship between appearance

690        and personality across the life span. *Personality and Social Psychology Bulletin, 24*,

691        736-749.

692    Zebrowitz, L. A., Franklin, R. G., Jr., & Boshyan, J. (2015). Face shape and behavior:

693        Implications of similarities in infants and adults. *Personality and Individual*

694        *Differences, 86*, 312-317.

695

696    **Figure Captions**

697

698    **Fig. 1.** Craniofacial landmarks used to calculate FWHR.

699    The skull width (the distance between the left and right zygions) is divided by the upper

700    facial height (the distance between the nasion and prosthion) to produce the FWHR. Figure

701    adapted from Weston et al. (2007).

702

703    **Supplementary Materials**

704

705    InformationOnDatasets.xlsx

706    These spreadsheets provide information on the populations included in the meta-analyses.

707

708    **S1 Fig.** Effect sizes for the 87 populations included in the skull meta-analysis.

709    The mean weighted effect size is highlighted in grey on the left. The eight outlying effect

710    sizes are labelled with red arrows.