

**Original citation:**

Warrington, N. M., Shevroja, E., Hemani, G., Hysi, P. G., Jiang, Y., Auton, A., Boer, C. G., Mangino, M., Wang, C. A., Kemp, J. P. [et al.](#) (2018) *Genome-wide association study identifies nine novel loci for 2D:4D finger ratio, a putative retrospective biomarker of testosterone exposure in utero*. Human Molecular Genetics, 27 (11). pp. 2025-2038. doi:[10.1093/hmg/ddy121](https://doi.org/10.1093/hmg/ddy121)

**Permanent WRAP URL:**

<http://wrap.warwick.ac.uk/100997>

**Copyright and reuse:**

The Warwick Research Archive Portal (WRAP) makes this work by researchers of the University of Warwick available open access under the following conditions. Copyright © and all moral rights to the version of the paper presented here belong to the individual author(s) and/or other copyright owners. To the extent reasonable and practicable the material made available in WRAP has been checked for eligibility before being made available.

Copies of full items can be used for personal research or study, educational, or not-for profit purposes without prior permission or charge. Provided that the authors, title and full bibliographic details are credited, a hyperlink and/or URL is given for the original metadata page and the content is not changed in any way.

**Publisher's statement:**

This is a pre-copyedited, author-produced PDF of an article accepted for publication in Human Molecular Genetics following peer review. The version of record Warrington, N. M., Shevroja, E., Hemani, G., Hysi, P. G., Jiang, Y., Auton, A., Boer, C. G., Mangino, M., Wang, C. A., Kemp, J. P. [et al.](#) (2018) *Genome-wide association study identifies nine novel loci for 2D:4D finger ratio, a putative retrospective biomarker of testosterone exposure in utero*. Human Molecular Genetics, 27 (11). pp. 2025-2038 is available online at: <https://doi.org/10.1093/hmg/ddy121>

**A note on versions:**

The version presented here may differ from the published version or, version of record, if you wish to cite this item you are advised to consult the publisher's version. Please see the 'permanent WRAP URL' above for details on accessing the published version and note that access may require a subscription.

For more information, please contact the WRAP Team at: [wrap@warwick.ac.uk](mailto:wrap@warwick.ac.uk)

**Title:**

Genome-wide association study identifies nine novel loci for 2D:4D finger ratio, a putative retrospective biomarker of testosterone exposure *in utero*

**Authors:**

Nicole M. Warrington<sup>1,2,3</sup>, Enisa Shevroja<sup>4,5</sup>, Gibran Hemani<sup>6,7</sup>, Pirro G. Hysi<sup>8</sup>, Yunxuan Jiang<sup>9</sup>, Adam Auton<sup>9</sup>, Cindy G. Boer<sup>5</sup>, Massimo Mangino<sup>8</sup>, Carol A. Wang<sup>3</sup>, John P. Kemp<sup>1,6,7</sup>, George McMahon<sup>6,7</sup>, Carolina Medina-Gomez<sup>4,5,10</sup>, Martha Hickey<sup>11</sup>, Katerina Trajanoska<sup>5,10</sup>, Dieter Wolke<sup>12</sup>, M. Arfan Ikram<sup>10</sup>, the 23andMe Research Team<sup>9</sup>, Grant W. Montgomery<sup>13</sup>, Janine F. Felix<sup>4,10,14</sup>, Margaret J. Wright<sup>13</sup>, David A. Mackey<sup>15</sup>, Vincent W. Jaddoe<sup>4,10,14</sup>, Nicholas G. Martin<sup>2</sup>, Joyce Y. Tung<sup>9</sup>, George Davey Smith<sup>6,7</sup>, Craig E. Pennell<sup>3</sup>, Tim D. Spector<sup>8</sup>, Joyce van Meurs<sup>4</sup>, Fernando Rivadeneira<sup>4,5,10</sup>, Sarah E. Medland<sup>2,¶</sup>, David M. Evans<sup>1,6,7,¶\*</sup>

**Affiliations:**

<sup>1</sup> The University of Queensland Diamantina Institute, The University of Queensland, Translational Research Institute, Brisbane, Queensland, 4102, Australia

<sup>2</sup> Queensland Institute of Medical Research, Brisbane, Queensland, 4006, Australia

<sup>3</sup> Division of Obstetrics and Gynaecology, The University of Western Australia, Perth, Western Australia, 6009, Australia

<sup>4</sup> The Generation R Study Group, Erasmus Medical Center, 3015, Rotterdam, The Netherlands

<sup>5</sup> Department of Internal Medicine, Erasmus Medical Center, 3015, Rotterdam, The Netherlands

<sup>6</sup> MRC Integrative Epidemiology Unit, University of Bristol, Bristol, BS8 2BN, England

<sup>7</sup> Population Health Sciences, University of Bristol, Bristol, BS8 2PS, England

<sup>8</sup> Department of Twin Research and Genetic Epidemiology, King's College London, London, SE1 7EH, England

<sup>9</sup> 23andMe, Inc., Mountain View, CA, 94061, USA

<sup>10</sup> Department of Epidemiology, Erasmus Medical Center, 3015, Rotterdam, Netherlands

<sup>11</sup> Department of Obstetrics and Gynaecology, The University of Melbourne and the Royal Women's Hospital, Parkville, Victoria, 3052, Australia

<sup>12</sup> Department of Psychology and Warwick Medical School, University of Warwick, Coventry CV47AL, England

<sup>13</sup> Queensland Brain Institute and Centre for Advanced Imaging, University of Queensland, Brisbane, Queensland, 4062, Australia

<sup>14</sup> Department of Pediatrics, Erasmus Medical Center, 3015, Rotterdam, The Netherlands

<sup>15</sup> Lions Eye Institute, Centre for Ophthalmology and Visual Science, University of Western Australia, Perth, Western Australia, 6009, Australia

<sup>¶</sup> These authors contributed equally to this work

\*Author for correspondence: David Evans, The University of Queensland Diamantina Institute, Level 7, 37 Kent Street, Translational Research Institute (TRI), Woolloongabba, Queensland 4102, Australia. Tel: +61-7-3443-7051; Fax: +61-7-3443-6966; Email: [d.evans1@uq.edu.au](mailto:d.evans1@uq.edu.au)

**Word count of abstract:** 219

**Word count of body (excluding tables/figures):**6,265

**Number of tables/figures:** 3/2

**Abstract:**

The ratio of the length of the index finger to that of the ring finger (2D:4D) is sexually dimorphic and is commonly used as a noninvasive biomarker of prenatal androgen exposure. Most association studies of 2D:4D ratio with a diverse range of sex-specific traits have typically involved small sample sizes and have been difficult to replicate, raising questions around the utility and precise meaning of the measure. In the largest genome-wide association meta-analysis of 2D:4D ratio to date (N=15,661, with replication N=75,821), we identified eleven loci (nine novel) explaining 3.8% of the variance in mean 2D:4D ratio. We also found weak evidence for association ( $\beta=0.06$ ;  $P=0.02$ ) between 2D:4D ratio and sensitivity to testosterone (length of the CAG microsatellite repeat in the androgen receptor gene) in females only. Furthermore, genetic variants associated with (adult) testosterone levels and/or sex hormone-binding globulin were not associated with 2D:4D ratio in our sample. Although we were unable to find strong evidence from our genetic study to support the hypothesis that 2D:4D ratio is a direct biomarker of prenatal exposure to androgens in healthy individuals, our findings do not explicitly exclude this possibility, and pathways involving testosterone may become apparent as the size of the discovery sample increases further. Our findings provide new insight into the underlying biology shaping 2D:4D variation in the general population.

## **Introduction:**

It has long been hypothesized that prenatal sex steroids, particularly testosterone, permanently modify the developing nervous system during critical periods of development, which in turn influences behavior in later life(1). Whilst animal models have largely supported this “Organizational hypothesis”(2), evidence from human studies has been much more limited, as accurately measuring prenatal testosterone exposure is extremely difficult. Based on several lines of indirect evidence, it has widely been hypothesized that the ratio of the length of the index finger to the length of the ring finger (2D:4D) is a marker of prenatal androgen exposure(3) and could therefore be used as a retrospective non-invasive biomarker of prenatal testosterone exposure. Males are believed to have greater prenatal testosterone exposure than females, and this is thought to determine the consistently observed lower ratio of the first (index) finger to the third (ring) finger in males (4, 5). This sex difference is relatively stable over time(3, 6, 7), and although there is variation in 2D:4D ratio across ethnic groups(8, 9), sexual dimorphism in the digit ratio is consistent across ethnicities(10). Based on the assumption that the digit ratio is a marker of prenatal testosterone exposure, associations have been reported between 2D:4D ratio and a broad range of sex-dependent behaviors and diseases including academic(11) and sporting performance(12, 13), social behaviors(14-16), fertility(17, 18), Alzheimer’s disease(19), metabolic syndrome indices(20) and autism spectrum disorder(21). Despite the extensive literature regarding 2D:4D ratio(22), most published studies have used small sample sizes (often fewer than 100 individuals) and have been difficult to replicate, raising questions around the utility and precise meaning of the measure(23).

Several twin studies have indicated that 2D:4D ratio is highly heritable ( $h^2$ : 50-80%)(24-28). Two relatively small genome-wide association studies (GWAS) of 2D:4D ratio have reported two loci influencing variation in the trait(29, 30). The minor allele (A) at rs314277 in *LIN28B* was associated with increased 2D:4D ratio, delayed menarche in females(31) and increased height(32). The minor allele (T) at the second locus, rs4902759 in *SMOC1*, was associated with decreased 2D:4D ratio.

Pedigree studies have shown that mutations in *SMOC1* are associated with Waardenburg anophthalmia (OMIM 206920), a syndrome that commonly includes abnormal digits(33-35). Additionally, the protein encoded by *SMOC1* has been shown to be up-regulated by androgens(36, 37) and down-regulated by estrogen(38), suggesting that *SMOC1* could be an intermediate between prenatal sex hormones and digit ratio(29).

Individuals with Complete Androgen Insensitivity Syndrome (CAIS) exhibit more feminine 2D:4D ratios, consistent with an effect of reduced prenatal testosterone exposure on digit ratio(39). These individuals have mutations in the androgen receptor gene (*AR*), located on the X chromosome, which codes for a receptor protein that facilitates physiological responses to androgens such as testosterone and dihydrotestosterone(40). *In vitro* studies suggest that the variable number of CAG repeats in exon 1 of *AR* is inversely related to the efficiency with which the receptor complex binds to DNA and influences transcription(41). Therefore, an association between the number of CAG repeats at the *AR* locus and 2D:4D ratio may indicate that sensitivity to androgens is a major driver of the individual differences in 2D:4D ratio seen in the normal population(40).

The aim of the present study was to investigate the genetic determinants of 2D:4D ratio by performing the largest GWAS meta-analysis (N=15,661, with replication N=75,821) to date. In addition, we aim to leverage genetics to scrutinize the evidence surrounding the hypothesis that 2D:4D ratio reflects prenatal androgen exposure. Specifically, we investigated whether there was association between repeat number in the androgen receptor (*XAR*) and 2D:4D ratio and we used our GWAS results to examine: 1) whether there was any association between 2D:4D ratio and genetic variants in pathways known to be linked with androgens; and 2) whether there was any relationship between 2D:4D ratio and genetic markers known to be related to (adult) levels of testosterone and/or sex hormone binding globulin (SHBG). We hypothesize that if 2D:4D ratio is truly influenced by levels of prenatal testosterone in utero, then it is reasonable to expect that genetic variants related to androgen

sensitivity (XAR) and/or serum levels of testosterone/SHBG might also show association with 2D:4D ratio.

## **Results:**

### *Study population*

Table 1 provides a summary of the studies included in the meta-analysis. Full details of the studies in the discovery meta-analysis and replication, including how the 2D:4D ratio was measured and the genotyping methods, are provided in the Supplementary Material. As expected, females had greater 2D:4D ratios than males in all cohorts. The mean and standard deviation of 2D:4D ratios measured on skeletal images (the Generation R Study and the Rotterdam Study) were lower than those of the other cohorts measured on photocopies of the hand, consistent with previous reports using this measure(42).

### *Genome-wide complex trait analysis of 2D:4D ratio of the left hand, right hand and mean of both hands*

Univariate genetic restricted maximum likelihood (GREML) analysis in ALSPAC revealed that common genome-wide variation explained a substantial proportion of the variance in 2D:4D ratio (left:  $h^2_{\text{SNP}}=0.299$ ,  $SE=0.071$ ,  $P=4.6 \times 10^{-6}$ ; right:  $h^2_{\text{SNP}}=0.360$ ,  $SE=0.071$ ,  $P=6.2 \times 10^{-8}$ ; mean:  $h^2_{\text{SNP}}=0.373$ ,  $SE=0.071$ ,  $P=1.4 \times 10^{-8}$ ). The genetic correlation between the left and right hand ratios was not different from 1 ( $r_g=0.918$ ,  $SE=0.074$ ,  $P=0.14$ ), indicating that the vast majority of SNPs contributing to variation in the ratio influence both hands. Therefore, in the main text we present the results from the left (typically non-dominant) hand ratio, due to the larger sample size achieved by including the Generation R Study, and include the right hand and mean ratio results in the Supplementary Material.

### *New genetic loci associated with 2D:4D ratio*

The meta-analysis of approximately 8.4 million 1000 genomes-imputed SNPs, including SNPs on the X chromosome, indicated that the lowest observed P-values for each of the three ratios deviated from

the expected null distribution (S1-Figure), whereas systematic inflation of the test statistics due to bias was negligible ( $\lambda_{\text{left}}=1.020$ ,  $\lambda_{\text{right}}=1.012$ ,  $\lambda_{\text{mean}}=1.014$ ). Further, no evidence of heterogeneity was detected between the discovery cohorts (S2-Figure). Eleven genomic loci reached genome-wide significance ( $P<5\times 10^{-8}$ ) in the discovery meta-analysis of 15,661 individuals for 2D:4D ratio (Figure 1 for Manhattan plots from the GWAS of the left hand [European only], S3-Figure for the left hand [multiethnic], right hand and average of both hands). Conditional and joint analysis in the genome-wide complex trait analysis (GCTA) software(43) identified two independent signals in the 16q12.1 locus, totaling 12 independent signals across the three 2D:4D ratio measurements. All 12 signals were replicated in 75,821 (52.8% male) research participants from 23andMe, Inc. (all  $P<0.004$ ; Table 2). Of the nine loci reaching genome-wide significance for the first time, six have not previously been described in the context of 2D:4D ratio, including: rs11581730 on chromosome 1q22; rs12474669 on chromosome 2q31.1; rs77640775 on chromosome 7p14.1; rs10790969 on chromosome 11q24.3; rs6499762 and rs1080014 on chromosome 16q12.1; and rs4799176 on chromosome 18q23. Two of the nine novel loci were reported at a suggestive significance level (but not genome-wide) in Medland *et al.*(30): SNPs in *LDAH* (previously known as *C2orf43*) on chromosome 2p24.1 and in *GLIS1* on chromosome 1p32.3. The locus on chromosome 2q31.1 is near the *HOXD* cluster of genes that are hypothesized to be required for growth and patterning of the digits, but this is the first time that convincing evidence for genetic association with 2D:4D ratio has been obtained. The remaining two loci included SNPs in *LIN28B*, identified previously by Medland *et al.*(30) and *SMOC1*, reported by Lawrance-Owen *et al.*(29). A summary of the meta-analysis results for the lead SNPs at each locus that reached genome-wide significance in the discovery sample are provided in Table 2, with the results from each study presented in S3-Table and regional plots in S4-Figure. The lead SNPs at the twelve replicated signals together explain 3.8% of the variance in 2D:4D ratio; this is equivalent to over half of the variance explained by sex in the Raine Study (5.1%).



We conducted analysis of the X chromosome to investigate whether there was evidence for 2D:4D ratio being partly an X-linked trait. No SNPs on the X chromosome reached genome-wide or suggestive significance (Figure 1). The *HOXA* gene cluster, along with the *HOXD* gene cluster, plays an important role in limb development, and were initially thought to be essential for the development of the 2D:4D ratio. Although the *HOXD* region was associated with 2D:4D ratio in our meta-analysis, there was no strong evidence for association between variants in the *HOXA* cluster and 2D:4D ratio (439 SNPs with all  $P > 0.001$ , S5-Figure).

As a secondary analysis, we also conducted sex-stratified analyses in each of the discovery cohorts and combined the results using a fixed-effects meta-analysis. The majority of the loci reaching genome-wide significance in the male or female only analyses were identified in the combined analysis (Miami plots in S6-Figure). One novel locus reached genome-wide significance in the female only analysis of the left hand 2D:4D ratio in the multi-ethnic meta-analysis. The top SNP in this region, rs10105686 (C allele [allele frequency = 0.79] Females:  $\beta = -0.334$ ,  $P = 2.42 \times 10^{-9}$ ; Males:  $\beta = -0.024$ ,  $P = 0.71$ ), is in *FGFR1* on chromosome 8. However, this locus only reached genome-wide significance for the left hand 2D:4D ratio in the multi-ethnic meta-analysis ( $P = 3.4 \times 10^{-7}$  in the female left hand European analysis), and would not be declared significant after correction for multiple testing given the large numbers of secondary analyses performed (i.e. secondary analyses involved analysis of males, females, left hand [both European and multi-ethnic analyses], right hand and average 2D:4D ratio, plus four sets of genome-wide sex heterogeneity analysis).

#### *Gene by sex interaction*

GREML analysis in ALSPAC showed no significant indication of gene by sex interaction (left:  $v_{gxe} = 0.000$ ,  $SE = 0.138$ ,  $P = 0.5$ ; right:  $v_{gxe} = 0.117$ ,  $SE = 0.144$ ,  $P = 0.21$ ; mean:  $v_{gxe} = 0.000$ ,  $SE = 0.140$ ,  $P = 0.5$ ). Consistent with this, only one locus on chromosome 9 reached genome-wide significance for difference in the magnitude of the regression coefficients between males and females for the average of both hands

2D:4D ratio (top SNP, rs16929125, A allele frequency = 0.91, heterogeneity  $P=1.17 \times 10^{-8}$ ; S7-Figure for Manhattan plots and S8-Figure for QQ plots). However, this locus only reached genome-wide significance for average 2D:4D ratio, and would not be declared significant after correction for multiple testing given the large number of secondary analyses performed.

#### *Gene prioritization, pathway and tissue analysis*

We used Data-driven Expression-Prioritized Integration for Complex Traits (DEPICT)(44) to identify the most likely causal gene at each locus and to investigate enriched pathways. DEPICT identified the nearest gene to the top associated signal to be the most likely causal gene in 10 of our 12 signals (S4-Table); *GLIS1*, *LDAH* (previously known as *C2orf43*), *OLA1*, *LIN28B*, *GLI3*, *FLI1*, *SMOC1*, *SALL1*, *TOX3* and *SALL3*. At the 2q31.1 locus, DEPICT prioritized two genes, *HODX11* and *HOXD12*, whilst at the 1q22 locus, five genes were prioritized: *EFNA1*, *DPM3*, *EFNA3*, *KRTCAP2* and *SLC50A1*. We will subsequently refer to this locus as *EFNA1*, which is the nearest gene to the top association signal.

When using the meta-analysis results from the average 2D:4D ratio of both hands, one gene set reached a false discovery rate (FDR)  $P<0.01$ , which mapped to the MSX1 PPI sub-network (S5-Table). The tissue enrichment analysis did not identify any tissues with a FDR  $P<0.01$  (S6-Table). Based on the expression data of 53 tissue types from the GTEx Consortium, four of the nearest genes to our lead SNPs showed high tissue expression in the testis or adrenal gland (*LDAH*, *LIN28B*, *SMOC1* and *C16orf97*; S9-Figure) relative to the other available tissues. Three of these four also showed expression in the brain (*LDAH*, *LIN28B* and *SMOC1*), in addition to two nearest genes, which showed high expression in the brain (*OLA1* and *SALL3*).

#### *Association between 2D:4D ratio and testosterone sensitivity*

We examined the association of 2D:4D ratio variation with the length of an established CAG repeat polymorphism in the AR gene, a proxy of testosterone sensitivity, in a meta-analysis of the ALSPAC

and QIMR cohorts. We found nominal evidence for a weak positive association between the number of CAG repeats in the *AR* gene on the X chromosome and mean 2D:4D ratio in females (mean of repeats:  $\beta=0.056$ ,  $P=0.02$ ; lower length repeat:  $\beta=0.047$ ,  $P=0.03$ ), but not in males (c.f. Table 3, see S7-Table for left and right 2D:4D ratio). The 91 SNPs in *AR* in the GWAS showed little evidence for association with 2D:4D ratio (minimum P-Value=0.03; S10-Figure).

#### *2D:4D associated variants and other traits*

Given the putative relationship between prenatal testosterone levels and 2D:4D ratio, we also examined the association between five published SNPs (rs12150660, rs5934505, rs10822186, rs10822184 and rs72829446) shown to influence testosterone levels(45, 46) using our 2D:4D ratio meta-analysis of the discovery cohorts; one other reported SNP was not tested, rs6258, as it was excluded from our meta-analysis as the minor allele frequency was <1%. Due to the high correlation between testosterone and its principal binding protein, sex hormone-binding globulin (SHBG), we also tested the association between the 13 published loci for SHBG(47, 48) and 2D:4D ratio. We observed three associations with left hand 2D:4D ratio at  $P<0.05$ : the C allele at rs1641537 (allele frequency = 0.87) was associated with increased 2D:4D ratio ( $\beta=0.111$ ,  $P=0.05$ ) and increased SHBG, the T allele at rs1573036 (allele frequency = 0.39) was associated with decreased 2D:4D ratio ( $\beta=-0.077$ ,  $P=0.02$ ) and increased SHBG and the T allele at rs72829446 (allele frequency = 0.11) was associated with increased 2D:4D ratio ( $\beta=0.123$ ,  $P=0.04$ ) and increased testosterone (Figure 2) (i.e. two out of the three associations were in the opposite direction to expected). We also failed to detect enrichment for association with 2D:4D ratio over all 18 testosterone and SHBG SNPs (Fisher's combined probability test  $P=0.10$ ). Additionally, there was no difference in the male and female effect sizes from the sex-stratified analyses across the 18 SNPs (Fisher's combined probability test for the heterogeneity P-value  $P=0.23$ ). Only one SNP, rs3779195, showed heterogeneity between males and females with the SHBG increasing T allele negatively associated with 2D:4D ratio in males and positively associated with 2D:4D ratio in females ( $\beta_{\text{male}}=-0.183$ ,  $P_{\text{male}}=0.01$ ;  $\beta_{\text{female}}=0.053$ ,  $P_{\text{female}}=0.39$ ;  $P_{\text{het}}=0.01$ ).

Given the previous observational epidemiological associations between 2D:4D ratio and sex-dependent behaviours, we used publicly available GWAS summary results for a variety of traits to estimate genetic correlations with 2D:4D ratio using linkage-disequilibrium (LD) score regression(49). We present the genetic correlation results in S8-Table, but advise caution in their interpretation, considering the recommendation by Bulik-Sullivan *et al*/(49) regarding conducting genetic correlation analysis for traits with heritability z-scores below 4, as the estimates tend to be noisy and less reliable (the z-scores for the left hand, right hand and average of both hands were 2.9, 3.0 and 3.3 respectively).

#### **Discussion:**

The 2D:4D ratio, a sexually dimorphic trait, has been extensively used in adults as a biomarker for prenatal androgen exposure. In the largest genetic association study of 2D:4D ratio to date, we identified nine novel loci for 2D:4D ratio, in addition to replicating two previously identified loci, *LIN28B* and *SMOC1*. These eleven loci explained 3.8% of the variance in mean 2D:4D ratio. After assessing association between 2D:4D ratio and a range of testosterone related traits, we found no conclusive evidence of the 2D:4D ratio constituting a marker of prenatal androgen exposure, although it is possible that pathways involving testosterone may become apparent as the size of our GWAS increases in the future. Yet, associations at distinct novel loci provide additional insight into the underlying biology shaping 2D:4D ratio variation.

The association signal on 1p32.3 spans *GLIS1*, which is expressed across several organs in the reproductive system, including the prostate, vagina, testis and cervix. Glis1, the protein encoded by *GLIS1*, is a Kruppel-like zinc finger protein that appears to have a critical role in controlling gene expression during specific stages of embryogenesis(50).

At 1q21-q22, the closest gene to the associated variants is *EFNA1*, which encodes a member of the ephrin family and has been implicated in mediating developmental events, notably in the nervous system. SNPs in the region of *EFNA1*, *DPM3* and *KRTCAP3* have previously been associated with prostate cancer risk(51) and  $\gamma$ -glutamyl transferase (GGT), an indicator of liver disease(52). Based on the results from GTEx, *EFNA1* is also mainly expressed in the liver, which is the most active site of lipid metabolism. A proxy for our lead 2D:4D ratio associated SNP at this locus, rs11264329, was associated with total and LDL cholesterol levels(53). Additionally, the top SNP identified at 2p24.1 is in *LDAH* (previously known as *C2orf43*); SNPs in *LDAH* have also been associated with prostate cancer risk(54, 55). The protein encoded by *LDAH* is involved in cholesterol mobilization(56). Testosterone, which is linked to prostate cancer risk, is created when luteinizing hormone (LH) triggers the testicular Leydig cells to convert cholesterol to testosterone. Therefore, these two loci could implicate cholesterol metabolism in steroidogenesis as a link to testosterone exerting a role on 2D:4D ratio variation, albeit requiring further investigation into the functional implications.

*OLA1*, which maps to the 2q31.1 locus, plays multiple roles in the regulation of cell proliferation and cell survival. Ding *et al.* show that mouse embryos lacking *OLA1* have delayed development leading to immature organs and stunted growth, which were frequently lethal prenatally(57). Their data suggests that there is a defect in cell proliferation due to a delay in cell cycle progression meaning that the mutant embryos appeared to undergo fewer proliferation cycles resulting in growth restriction.

SNPs in the 7p14.1 region map within *GLI3*. The gene encodes a zinc finger transcription factor that functions in the hedgehog signal transduction pathway. SNPs in this region have also been associated with facial morphology, namely nose wing breadth(58), and implicated in several Mendelian disorders which are characterized by craniofacial and limb abnormalities. Specifically, there are several disorders and conditions where polydactyly is a feature(59), including Greig cephalopolysyndactyly syndrome (GCPS), and Pallister-Hall syndrome (PHS). There is some evidence for brachydactyly

(shortened digits) in a mouse null *Gli3* mutant developed by Sheth *et al.*(60), and in patients with PHS(59).

The association signal arising from the 11q24.1-q24.3 locus is intergenic between *FLI1* and *ETS1*. Through GWAS, SNPs in the *FLI1* gene have been shown to be associated with height, with similar effects in both males and females(61). In addition, SNPs in *ETS1* have been shown to be associated with rheumatoid arthritis(62) and celiac disease(63) in European populations and with systemic lupus erythematosus(64-66) in Chinese populations; all of these diseases have a higher prevalence in females. This gene encodes the protein Ets1, which is expressed in a variety of tissues throughout the development of an embryo and plays a role in pituitary hormone secretion(67). In mice, the ETS factor family defines Shh spatial expression in limb buds and alterations define pathogenetic mechanism leading to preaxial polydactyly(68).

The GWAS signal on chromosome 16q12.1 maps in the vicinity of the *SALL1*, *TOX3* and *C16orf97* genes. Not much is known about the function of *C16orf97*. However, *SALL1* was identified by DEPICT as being the most likely causal gene for the rs6499762 association; mutations in *SALL1* cause Townes-Brocks syndrome(69), a condition characterized by hand malformations, abnormally shaped ears and anal atresia, among other genital malformations(70). Additionally, *SALL3* on chromosome 18q23, is also part of the human *spalt*-like gene family, which is associated with syndromic forms presenting with skeletal abnormalities. Kohlhasse *et al.*(71) characterized this gene and implicated it in the 18q deletion syndrome, which results in mental and growth retardation, developmental delay, hearing loss, and facial and limb abnormalities including tapered fingers(72). Altogether, several links between 2D:4D ratio and testosterone metabolism can be derived from these associations, involving hormonal pathways and the process of sexual differentiation during early development. Further, SNPs in the *SALL3* region have also been associated with prostate cancer(51), which may suggest additional links between 2D:4D ratio and testosterone.

One gene set was identified as being associated with 2D:4D ratio, the MSX1 PPI subnetwork. An Msx1-interacting network of transcription factors has been shown to operate during early tooth development(73). Therefore, the identification of this gene set may be highlighting a network that is involved in several areas of development.

Interestingly, we didn't find strong evidence of association between variants within the *HOXA* gene cluster and 2D:4D ratio. This lack of association does not preclude variation in distal enhancers acting through effects on the expression of *HOXA* cluster genes nor that variants of smaller effect acting from within the cluster itself (with an alpha of  $5 \times 10^{-8}$  we had 80% power to detect a genetic variant that explained approximately 0.28% of the variance in the left hand 2D:4D ratio in Europeans only (N=14,382)). We did, however, find an association involving a variant in *HOXD12* (the most strongly associated SNP was rs847158,  $P=9.58 \times 10^{-11}$ , which replicated in the 23andMe dataset). The exact role of *HOXD12* has not yet been determined; however, the homeobox family of genes plays an important role in morphogenesis and is particularly relevant in the development of the limbs and genitals(74).

#### *2D:4D as a marker of testosterone exposure*

2D:4D ratio has been used extensively in adults as a biomarker for prenatal androgen exposure. However, whether the digit ratio reliably reflects prenatal androgen exposure has not been convincingly demonstrated. Most of the data linking 2D:4D ratio with prenatal androgen exposure is based on preclinical or indirect evidence, including studies that indicate that 2D:4D ratio is fixed early in gestation and is associated with adult levels of circulating testosterone(3, 75). The most direct test of this hypothesis to date was performed by Lutchmaya *et al.* who showed that testosterone levels in amniotic fluid from the second trimester of pregnancy were not associated with 2D:4D ratio at two years of age(76). However, the authors did find that an increased ratio of testosterone to estradiol was associated with a lower (or more male like) 2D:4D ratio, suggesting that the relationship between

digit ratio and prenatal hormones may be more complicated and not only reflect testosterone levels(76, 77).

In the present study we attempted to use genetic evidence to find support for the testosterone biomarker hypothesis. Our rationale was that if prenatal testosterone affects 2D:4D ratio, then it is logical that polymorphisms in genes related to androgen sensitivity (e.g. in *XAR*) and/or SNPs robustly associated with androgen levels/levels of SHBG, should also be related to 2D:4D ratio. Whilst we did detect some evidence of a positive association between the number of CAG repeats in *AR* and 2D:4D ratio in females, we note that the small effect size would not be significant after adjusting for the multiple statistical tests we performed. Power calculations suggest that our combined sample of N=5,826 individuals in the *XAR* meta-analysis was well powered (~78%) to detect a locus responsible for 0.001% of the phenotypic variance in 2D:4D ratio (one tailed  $\alpha=0.05$ ). In comparison, all of our genome-wide significant SNPs explained >0.001% of the variance in 2D:4D ratio (most explained much more variance than this). This suggests that if genetic variation in *XAR* does contribute to variation in 2D:4D ratio through, for example, sensitivity to testosterone, its effect is likely to be small relative to other sources of genetic variation. Two recent smaller meta-analyses also failed to find an association between length of the repeat in *XAR* and 2D:4D ratio(78, 79).

Likewise, using SNPs that are associated with testosterone and SHBG we were unable to detect any enrichment for association with 2D:4D ratio. We were also unable to identify any genetic correlation between 2D:4D ratio and a range of traits and diseases previously implicated with 2D:4D ratio variation. This indicates that the previously identified observational associations may not be driven by known genetic loci that are shared between the traits, although we acknowledge the power of analysis was low and confidence intervals around our estimates were large.



372 Whilst we were unable to find any convincing evidence that sensitivity to/levels of androgens is a  
373 major driver of the individual differences in 2D:4D ratio seen in the normal population, there are  
374 several key assumptions underlying the use of genetic variation to investigate the link between  
375 prenatal androgen exposure and 2D:4D ratio. First, investigating the association between the number  
376 of CAG repeats in AR and 2D:4D ratio relies on the assumption that CAG length reflects androgen  
377 sensitivity. There is fairly good evidence for this, at least *in vitro* as derived from the Chamberlain *et*  
378 *al.* functional study showing a linear relationship between increased CAG length and decreased  
379 transactivation function(41). Secondly, we assume that the genotyping of the CAG repeat is accurate.  
380 Although there was some discordance in the replicate genotyping in ALSPAC, the majority of  
381 discrepancies were only one CAG repeat different between the replicates. Our simulations presented  
382 in the Supplementary Material indicate that this degree of measurement error had little influence on  
383 the power of our association analysis. Thirdly, we assume that the SNPs associated with adult levels  
384 of testosterone/SHBG also reflect testosterone/SHBG levels prenatally and that the effect of these  
385 SNPs are similar in males and females. The testosterone associated SNPs were identified in GWAS of  
386 adult men only and one of the two GWAS for SHBG, identifying only one novel locus, was in post-  
387 menopausal women only. However, Coviello *et al.* conducted sex-stratified analyses and identified  
388 only one locus with significant heterogeneity on SHBG between males and females (48). They also  
389 showed that the SHBG SNPs identified explained ~15.6% of the variation in SHBG in men and ~8.4% of  
390 the variation in women. This indicates that although the SNPs have a greater overall effect in males,  
391 they are still likely to be associated with androgens in females. It is too difficult to measure androgen  
392 levels in the fetus so we are unable to confirm that these SNPs are also associated with androgen  
393 levels prenatally. Finally, although none of our novel loci showed direct evidence of being related to  
394 pathways involving testosterone, it does not preclude the very real possibility that testosterone  
395 influences 2D:4D ratio by downregulating or upregulating the expression of genes involved in its  
396 determination. Indeed, in GWAS of height (80)and BMI (81)the role of expected hormone pathways

only appeared when the studies were sufficiently powered to find far greater number of genome-wide significant hits than in the present study.

In conclusion, we have conducted the largest GWAS of 2D:4D ratio to date and identified nine novel loci robustly associated with 2D:4D ratio in Europeans, bringing the total number of robustly associated loci to eleven. We were unable to find strong evidence from our genetic study to support the hypothesis that 2D:4D ratio is a direct biomarker of prenatal exposure to androgens in healthy individuals, although our findings do not explicitly exclude this possibility, and pathways involving testosterone may become apparent as the size of the discovery sample increases further. Our findings provide new insight into the underlying biology shaping 2D:4D variation in the general population.

## **Materials and Methods:**

### *Participants*

We drew on data from six cohorts for the discovery genome-wide meta-analysis and association with the CAG repeat in *AR* including the Avon Longitudinal Study of Parents and Children (ALSPAC), the Generation R Study, the Rotterdam Study, the Western Australian Pregnancy Cohort (Raine) Study, TwinsUK and the Queensland Institute of Medical Research (QIMR) sample, which was drawn from the Brisbane Adolescent Twin Study (BATS; also known as the Brisbane Longitudinal Twin Study (BLTS)). The 23andMe cohort was used for replication of the genome-wide significant findings. Details of each of these studies, including how the 2D:4D ratio was measured and the genotyping methods, are provided in the Supplementary Material.

### *Ethics Statement*

All cohorts in the discovery meta-analysis or replication obtained ethical approval from their local ethics review boards; ALSAPC from the ALSPAC Law and Ethics Committee and the Local Research Ethics Committees, The Generation R Study from the Medical Ethics Committee of the Erasmus Medical Center in Rotterdam, QIMR from the QIMR Human Research Ethics Committee, the Raine study from the King Edward Memorial Hospital and Princess Margaret Hospital for Children Human Research Ethics Committees, the Rotterdam Study from the Medical Ethics Committee of the Erasmus Medical Center in Rotterdam, TwinsUK from the Guy's and St Thomas' (GSTT) Ethics Committee and research participants from 23andMe provided informed consent and participated in the research online, under a protocol approved by the external AAHRPP-accredited IRB, Ethical & Independent Review Services (E&I Review).

## *Statistical Analysis*

The 2D:4D ratio was calculated as the length of the second digit divided by the length of the fourth digit, multiplied by 100 so as to avoid computational difficulties due to the low variance of the trait. In all studies, the measure was normally distributed, so no further transformation was applied.

### Genome-wide complex trait analysis of 2D:4D ratio:

To estimate the proportion of additive genetic variance in 2D:4D ratio explained by directly genotyped SNPs we conducted a univariate genetic restricted maximum likelihood (GREML) analysis using the genome-wide complex trait analysis (GCTA) software(82) in >4,900 individuals from ALSPAC. Sex was included as a fixed effect in the model. Bivariate GREML analysis(83) was used to estimate the genetic correlation between the left and right hand 2D:4D ratio, which will indicate whether the same genetic variants contribute to variation in the ratio of each hand. Additionally, a gene by sex analysis was conducted using the gene by environment test, to indicate whether the SNPs associated with 2D:4D ratio differed between males and females.

### Genome-wide association analysis - Discovery:

Genome-wide association analysis using imputed dosages to account for uncertainty in the imputation was performed using linear regression in each cohort, adjusting for sex. In addition, the QIMR and TwinsUK cohorts accounted for zygosity and relatedness. The Generation R Study (European subset), the Rotterdam Study and Raine adjusted for four principal components for population stratification. A sensitivity analysis to maximize power was conducted by including all individuals of the Generation R Study with adjustment for 20 principal components to account for the multi-ethnic sample as performed earlier(84-88). Results presented in the main text are derived from the meta-analysis including the European subset only, with the multi-ethnic analysis presented in the Supplementary Material. SNPs were tested for association with left 2D:4D ratio, right 2D:4D ratio (all cohorts excluding the Generation R Study) and the mean of the left and right hand 2D:4D ratios (all cohorts excluding

the Generation R Study). Results were combined using fixed-effects inverse-variance weighted meta-analysis in METAL(89), adjusting for genomic control. Within each study, SNPs with a MAF < 1%, an INFO score < 0.4 or R2 for imputation quality < 0.3 were excluded from the meta-analysis and SNPs that were reported in less than 50% of the total sample size were excluded from further follow-up.

We also tested whether the regression coefficients differed between males and females. We carried out genome-wide association analysis in males and females separately in each of the discovery cohorts and, as with the main analysis, we excluded variants with a MAF <1% and poor imputation quality (INFO score <0.4 or R2 for imputation quality < 0.3). We performed a fixed-effects inverse-variance weighted meta-analysis of each sex in METAL(89), adjusting for genomic control. Finally, we excluded variants that were reported less than 50% of the male and female sample sizes, and performed a chi-square test of heterogeneity between the meta-analyzed male and female effects in METAL(89) to test for the difference between the effect sizes in males and females and produce an overall level of significance.

#### Conditional and joint association analysis:

We performed approximate conditional and joint SNP association analysis using the GCTA software(43), which utilizes meta-analysis summary statistics and linkage disequilibrium (LD) structure from a reference sample. We used this approach to identify additional signals in regions of association, using a subset of 15,000 UK Biobank(90) individuals as the reference sample to approximate LD patterns. The selected subset of the UK Biobank individuals were of European descent and unrelated to anyone else in the subset.

#### Genome-wide association analysis - Replication:

SNPs that reached genome-wide significance ( $P < 5 \times 10^{-8}$  for the left, right or average 2D:4D ratio) in the discovery meta-analysis were replicated in the 23andMe dataset. If the imputed SNP passed quality

control or the genotyped SNP was unavailable, then the imputed SNP was used for analysis, otherwise the genotyped SNP was used. Analysis of each SNP was performed using linear regression, adjusting for sex, age, the first five principal components for population stratification and genotyping platform. Results were adjusted for a genomic control inflation factor of  $\lambda=1.074$ .

#### Variance explained:

The variance explained by each SNP was calculated using the effect size from the discovery meta-analysis (beta,  $\beta$ ), the minor and major allele frequencies ( $p$  and  $q$  respectively) and the variance of the 2D:4D ratio ( $Var(Y)$ ) using the following formula:

$$VarExp = 2pq \frac{\beta^2}{Var(Y)}$$

We used the median standard deviation in 2D:4D ratio across the cohorts, which was 3.40 for the left hand, 3.35 for the right hand and 3.06 for the average of both hands to calculate the phenotypic variance in this formula. Under the assumption that all SNPs independently contribute to 2D:4D ratio, we computed the total variance explained by the lead SNPs at the genome-wide significant loci as the sum of the single-SNP explained variances.

#### Gene prioritization, gene set and tissue/cell type enrichment analysis:

We conducted three analyses implemented in DEPICT (44) to establish the functional connections with our lead signals. First, we prioritized genes which are most likely to be causal for 2D:4D ratio by correlating the reconstituted gene set membership of each gene nearby the associated signal to genes from other associated loci and adjusting for potential sources of bias such as gene size. Second, we performed a gene set enrichment analysis, which tests if the genes in the associated loci are enriched in the reconstituted gene-sets. Third, we analysed expression enrichment across particular tissues or cell types, by testing whether genes associated with 2D:4D ratio loci were seen highly expressed in any of the 209 Medical Subject Heading (MeSH) annotations using data from 37,427 expression arrays.

In all three analyses, we used false discovery rate (FDR) to adjust for multiple testing, with an FDR P-value < 0.01 defined as significant.

The DEPICT analyses were based on independent lead SNPs ( $r^2 < 0.1$ , European populations 1000 genomes reference panel) with P-values below the genome-wide significant threshold ( $P < 5 \times 10^{-8}$ ). Gene-set enrichment was further grouped into 'meta gene sets' by similarity clustering, as previously described (44).

Additionally, we investigated the functionality of the genes closest to the lead SNPs identified in these analyses using expression data on 53 tissue types from the Genotype-Tissue Expression (GTEx) consortium(91)

#### Analysis of the CAG repeat polymorphism in AR:

Information on the CAG repeat polymorphism in AR was available in the ALSPAC and QIMR cohorts (see Supplementary Material for genotyping information). In ALSPAC, we performed linear regression of 2D:4D ratio (left, right and mean of the left and right) on length of CAG repeat (in females either average repeat length, the highest length repeat or the lower length repeat). The QIMR analyses were conducted using full information maximum likelihood structural equation models in openMx(92) which explicitly accounted for relatedness and zygosity while estimating the linear effect of the CAG repeat on 2D:4D ratio. We performed analyses including all participants (with sex as a covariate), females separately and males separately. A fixed-effects inverse-variance weighted meta-analysis was used to combine the results from the two cohorts using the rmeta package in R (version 3.0.0)(93). A one-tailed hypothesis was used to test whether there was a positive association between the number of CAG repeats and 2D:4D ratio.

#### Genetic correlation with associated traits:

We used LD score regression, which has been described in detail elsewhere<sup>(49)</sup>, to calculate the genetic correlation between 2D:4D ratio and a range of traits and diseases it has been associated with in observational studies. Note, we conducted this analysis using the European only meta-analysis results as LD score regression cannot accommodate LD variation between diverse populations. Briefly, the LD Score is a measure of how much genetic variation each SNP tags; so if a SNP has a high LD Score then it is in high LD with many nearby SNPs. SNPs with high LD Scores are more likely to contain more true signals and hence provide more chance for overlap with genuine signals between GWAS. The method uses summary statistics from the GWAS meta-analyses of 2D:4D ratio and the traits of interest, calculates the cross-product of test statistics at each SNP, and then regresses the cross-product on the LD Score. The slope of the regression is a function of the genetic correlation between traits. If there is overlap between the samples used in each of the meta-analyses (or cryptic relatedness between samples) it will only affect the intercept of the regression, and will not bias the estimate of the genetic covariance.

Summary statistics from the GWAS meta-analysis for traits and diseases of interest were downloaded from the relevant consortium website (see S8-Table for references). The summary statistics files were reformatted for LD Score regression analysis using the `munge_sumstats.py` python script provided on the developer's website (<https://github.com/bulik/ldsc>); we filtered the summary statistics to the subset of HapMap3 SNPs, as advised by the developers, to ensure that no bias was introduced due to poor imputation quality. Where the sample size for each SNP was included in the results file this was flagged using `--N-col`; if no sample size was available then the maximum sample size reported in the reference for the GWAS meta-analysis was used (i.e the summary statistics for each SNP was assumed to have been estimated using the same sample size). SNPs were excluded if the minor allele frequency was  $<0.01$ , the strand was ambiguous, the rs number was duplicated or they had a sample size less than 60% of the total sample size available. Once all the files were reformatted, we used the `ldsc.py`



python script, also on the developer's website, to calculate the genetic correlation between 2D:4D ratio and each of the traits and diseases. The European LD Score files that were calculated from the 1000 Genomes reference panel and provided by the developers were used for the analysis.

#### **Acknowledgements and Funding:**

N.M.W is supported by a National Health and Medical Research Council Early Career Fellowship (APP1104818). E.S. is supported by the European Commission within the framework of the Erasmus-Western Balkans (ERAWEB). M.H is supported by a National Health and Medical Research Council Practitioner Fellowship (APP1058935). F.R. is supported by the Netherlands Scientific Organization (NWO) and ZonMW Project number: NWO/ZONMW-VIDI-016-136-367. S.E.M is supported by a National Health and Medical Research Council Senior Research Fellowship (APP1103623). D.M.E is supported by an Australian Research Council Future Fellowship (FT130101709) and an MRC programme grant (MC\_UU\_12013/4).

ALSPAC: We are extremely grateful to all the families who took part in this study, the midwives for their help in recruiting them, and the whole ALSPAC team, which includes interviewers, computer and laboratory technicians, clerical workers, research scientists, volunteers, managers, receptionists and nurses. GWAS data was generated by Sample Logistics and Genotyping Facilities at the Wellcome Trust Sanger Institute and LabCorp (Laboratory Corporation of America) using support from 23andMe. The UK Medical Research Council and the Wellcome Trust (grant reference: 102215/2/13/2) and the University of Bristol provide core support for ALSPAC. This work is supported by a Medical Research Council program grant (grant reference:MC\_UU\_12013/4 to D.M.E). The androgen receptor CAG repeat data was generated with funding from the Medical Research Council (grant reference: G0500953 to Barbara Maughan). This publication is the work of the authors, and D.M.E will serve as guarantor for the contents of this paper.

Generation R Study: The Generation R Study is conducted by the Erasmus Medical Center in close collaboration with the School of Law and Faculty of Social Sciences of the Erasmus University Rotterdam, the Municipal Health Service Rotterdam area, Rotterdam, the Rotterdam Homecare Foundation, Rotterdam and the Stichting Trombosedienst & Artsenlaboratorium Rijnmond [STAR-MDC], Rotterdam. We gratefully acknowledge the contribution of children and parents, general practitioners, hospitals, midwives and pharmacies in Rotterdam. The generation and management of GWAS genotype data for the Generation R Study was done at the Genetic Laboratory of the Department of Internal Medicine, Erasmus MC, the Netherlands. We would like to thank Karol Estrada, Dr. Tobias A. Knoch, Anis Abuseiris, Luc V. de Zeeuw, and Rob de Graaf, for their help in creating GRIMP, BigGRID, MediGRID, and Services@MediGRID/D-Grid, [funded by the German Bundesministerium fuer Forschung und Technology; grants 01 AK 803 A-H, 01 IG 07015 G] for access to their grid computing resources. We thank Pascal Arp, Mila Jhamai, Marijn Verkerk, Manoushka Ganesh, Lizbeth Herrera and Marjolein Peters for their help in creating, managing and QC of the GWAS database.

The general design of Generation R Study is made possible by financial support from the Erasmus Medical Center, Rotterdam, the Erasmus University Rotterdam, the Netherlands Organization for Health Research and Development (ZonMw), the Netherlands Organisation for Scientific Research (NWO), the Ministry of Health, Welfare and Sport and the Ministry of Youth and Families. The musculoskeletal research of the Generation R Study is partly supported by the European Commission grant HEALTH-F2-2008-201865-GEFOS. Additionally, the Netherlands Organization for Health Research and Development supported authors of this manuscript (ZonMw 907.00303, ZonMw 916.10159, ZonMw VIDI 016.136.361 to V.W.V.J., and ZonMw VIDI 016.136.367 to F.R.). This project also received funding from the European Union's Horizon 2020 research and innovation programme under grant agreements No 633595 (DynaHEALTH) and No 733206 (LIFECYCLE), and from the European Research Council (ERC Consolidator Grant, ERC-2014-CoG-648916 to V.W.V.J.).

QIMR: We thank the Brisbane twins and siblings for their participation; Marlene Grace, Ann Eldridge and Natalie Garden for sample collection; Kerrie McAloney for study co-ordination; Harry Beeby, Daniel Park, and David Smyth for IT support, Anjali Henders and the Molecular Genetics Laboratory for DNA sample preparation, and Scott Gordon for genotyping QC. The QIMR studies were supported by funding from the Australian National Health and Medical Research Council (grant numbers: 241944, 339462, 389927, 389875, 389891, 389892, 389938, 443036, 442915, 442981, 496739, 552485, and 552498, and most recently 1049894) and the Australian Research Council (grant numbers: A7960034, A79906588, A79801419, DP0212016, and DP0343921).

The Rotterdam Study: The generation and management of GWAS genotype data for the Rotterdam Study (RS I, RS II, RS III) was executed by the Human Genotyping Facility of the Genetic Laboratory of the Department of Internal Medicine, Erasmus MC, Rotterdam, The Netherlands. The GWAS datasets are supported by the Netherlands Organisation of Scientific Research NWO Investments (nr. 175.010.2005.011, 911-03-012), the Genetic Laboratory of the Department of Internal Medicine, Erasmus MC, the Research Institute for Diseases in the Elderly (014-93-015; RIDE2), the Netherlands Genomics Initiative (NGI)/Netherlands Organisation for Scientific Research (NWO) Netherlands Consortium for Healthy Aging (NCHA), project nr. 050-060-810. We thank Pascal Arp, Mila Jhamai, Marijn Verkerk, Lizbeth Herrera and Marjolein Peters, MSc, and Carolina Medina-Gomez, MSc, for their help in creating the GWAS database, and Karol Estrada, PhD, Yurii Aulchenko, PhD, and Carolina Medina-Gomez, PhD, for the creation and analysis of imputed data. We would like to thank Dr. Karol Estrada, Dr. Fernando Rivadeneira, Dr. Tobias A. Knoch, Marijn Verkerk, Anis Abuseiris, Dr. Linda Boer and Rob de Graaf (Erasmus MC Rotterdam, The Netherlands), for their help in creating and maintaining GRIMP. Dr. Fernando Rivadeneira received an additional grant from the Netherlands Organization for Health Research and Development ZonMw VIDI 016.136.367. The Rotterdam Study is funded by Erasmus Medical Center and Erasmus University, Rotterdam, Netherlands Organization for the Health Research and Development (ZonMw), the Research Institute for Diseases in the Elderly

(RIDE), the Ministry of Education, Culture and Science, the Ministry for Health, Welfare and Sports, the European Commission (DG XII), and the Municipality of Rotterdam. The authors are very grateful to the study participants, the staff from the Rotterdam Study (particularly L. Buist and J.H. van den Boogert) and the participating general practitioners and pharmacists.

Twins UK: The study was funded by the Wellcome Trust (Ref: 105022/Z/14/Z); European Community's Seventh Framework Programme (FP7/2007-2013). The study also receives support from the National Institute for Health Research (NIHR) - funded BioResource, Clinical Research Facility and Biomedical Research Centre based at Guy's and St Thomas' NHS Foundation Trust in partnership with King's College London. SNP Genotyping was performed by The Wellcome Trust Sanger Institute and National Eye Institute via NIH/CIDR

Raine: The authors are grateful to the Raine Study participants, their families, and to the Raine Study research staff for cohort coordination and data collection. The authors gratefully acknowledge the assistance of the Western Australian DNA Bank (National Health and Medical Research Council of Australia National Enabling Facility). The following Institutions provide funding for Core Management of the Raine Study: The University of Western Australia (UWA), Raine Medical Research Foundation, UWA Faculty of Medicine, Dentistry and Health Sciences, The Telethon Institute for Child Health Research, Curtin University, Edith Cowan University and Women and Infants Research Foundation. This study was supported by project grants from the National Health and Medical Research Council of Australia (Grant numbers: 403981, 003209 and 1021105; <http://www.nhmrc.gov.au/>) and the Canadian Institutes of Health Research (Grant number: MOP-82893; <http://www.cihr-irsc.gc.ca/e/193.html>). This work was supported by resources provided by the Pawsey Supercomputing Centre with funding from the Australian Government and the Government of Western Australia. The funders had no role in study design, data collection and analysis, decision to publish, or preparation of the manuscript.

664 23andMe: We thank the 23andMe research participants and employees for their contributions to this  
665 study. We thank the following members of the 23andMe Research Team: Michelle Agee, Babak  
666 Alipanahi, Robert K. Bell, Katarzyna Bryc, Sarah L. Elson, Pierre Fontanillas, Nicholas A. Furlotte, David  
667 A. Hinds, Bethann S. Hromatka, Karen E. Huber, Aaron Kleinman, Nadia K. Litterman, Matthew H.  
668 McIntyre, Joanna L. Mountain, Elizabeth S. Noblin, Carrie A.M. Northover, Steven J. Pitts, J. Fah  
669 Sathirapongsasuti, Olga V. Sazonova, Janie F. Shelton, Suyash Shringarpure, Chao Tian, Vladimir Vacic,  
670 and Catherine H. Wilson.  
671

672 **References:**

- 673 1 Phoenix, C.H., Goy, R.W., Gerall, A.A. and Young, W.C. (1959) Organizing action of prenatally  
674 administered testosterone propionate on the tissues mediating mating behavior in the female  
675 guinea pig. *Endocrinology*, **65**, 369-382.
- 676 2 Arnold, A.P. and Breedlove, S.M. (1985) Organizational and activational effects of sex  
677 steroids on brain and behavior: a reanalysis. *Hormones and behavior*, **19**, 469-498.
- 678 3 Manning, J.T., Scutt, D., Wilson, J. and Lewis-Jones, D.I. (1998) The ratio of 2nd to 4th digit  
679 length: a predictor of sperm numbers and concentrations of testosterone, luteinizing hormone and  
680 oestrogen. *Human reproduction (Oxford, England)*, **13**, 3000-3004.
- 681 4 Ecker, A. (1875) Einige Bemerkungen u"ber einen schwankenden Charakter in der Hand des  
682 Menschen. *Archiv fu"r Anthropologie*, **8**, 67-74.
- 683 5 Wilson, G.D. (1983) Finger-length as an index of assertiveness in women. *Personality and*  
684 *Individual Differences*, **4**, 111-112.
- 685 6 McIntyre, M.H., Ellison, P.T., Lieberman, D.E., Demerath, E. and Towne, B. (2005) The  
686 development of sex differences in digital formula from infancy in the Fels Longitudinal Study.  
687 *Proceedings. Biological sciences / The Royal Society*, **272**, 1473-1479.
- 688 7 Trivers, R., Manning, J. and Jacobson, A. (2006) A longitudinal study of digit ratio (2D:4D) and  
689 other finger ratios in Jamaican children. *Hormones and behavior*, **49**, 150-156.
- 690 8 Manning, J.T., Henzi, P., Venkatramana, P., Martin, S. and Singh, D. (2003) Second to fourth  
691 digit ratio: ethnic differences and family size in English, Indian and South African populations. *Annals*  
692 *of human biology*, **30**, 579-588.
- 693 9 Loehlin, J.C., McFadden, D., Medland, S.E. and Martin, N.G. (2006) Population differences in  
694 finger-length ratios: ethnicity or latitude? *Archives of sexual behavior*, **35**, 739-742.
- 695 10 Manning, J.T., Barley, L., Walton, J., Lewis-Jones, D.I., Trivers, R.L., Singh, D., Thornhill, R.,  
696 Rohde, P., Bereczkei, T., Henzi, P. *et al.* (2000) The 2nd:4th digit ratio, sexual dimorphism, population  
697 differences, and reproductive success. evidence for sexually antagonistic genes? *Evolution and*  
698 *human behavior : official journal of the Human Behavior and Evolution Society*, **21**, 163-183.
- 699 11 Coco, M., Perciavalle, V., Maci, T., Nicoletti, F., Di Corrado, D. and Perciavalle, V. (2011) The  
700 second-to-fourth digit ratio correlates with the rate of academic performance in medical school  
701 students. *Molecular medicine reports*, **4**, 471-476.
- 702 12 Hull, M.J., Schranz, N.K., Manning, J.T. and Tomkinson, G.R. (2014) Relationships between  
703 digit ratio (2D:4D) and female competitive rowing performance. *American journal of human biology :*  
704 *the official journal of the Human Biology Council*, in press.
- 705 13 Trivers, R., Hopp, R. and Manning, J. (2013) A longitudinal study of digit ratio (2D:4D) and its  
706 relationships with adult running speed in Jamaicans. *Human biology*, **85**, 623-626.
- 707 14 Kim, Y., Kim, K. and Kim, T.H. (2014) Domain Specific Relationships of 2D:4D Digit Ratio in  
708 Risk Perception and Risk Behavior. *The Journal of general psychology*, **141**, 373-392.
- 709 15 Madison, G., Aasa, U., Wallert, J. and Woodley, M.A. (2014) Feminist activist women are  
710 masculinized in terms of digit-ratio and social dominance: a possible explanation for the feminist  
711 paradox. *Frontiers in psychology*, **5**, 1011.
- 712 16 Lam, D. and Ozorio, B. (2014) An Exploratory Study of the Relationship Between Digit Ratio,  
713 Illusion of Control, and Risk-Taking Behavior Among Chinese College Students. *Journal of gambling*  
714 *studies / co-sponsored by the National Council on Problem Gambling and Institute for the Study of*  
715 *Gambling and Commercial Gaming*, in press.
- 716 17 Oh, J.K., Kim, K.T., Yoon, S.J., Kim, S.W. and Kim, T.B. (2014) Second to fourth digit ratio: a  
717 predictor of adult testicular volume. *Andrology*, **2**, 862-867.
- 718 18 Klimek, M., Galbarczyk, A., Nenko, I., Alvarado, L.C. and Jasienska, G. (2014) Digit ratio  
719 (2D:4D) as an indicator of body size, testosterone concentration and number of children in human  
720 males. *Annals of human biology*, **41**, 518-523.

721 19 Vladeanu, M., Giuffrida, O. and Bourne, V.J. (2014) Prenatal sex hormone exposure and risk  
722 of Alzheimer disease: a pilot study using the 2D:4D digit length ratio. *Cognitive and behavioral*  
723 *neurology : official journal of the Society for Behavioral and Cognitive Neurology*, **27**, 102-106.

724 20 Oyeyemi, B.F., Iyiola, O.A., Oyeyemi, A.W., Oricha, K.A., Anifowoshe, A.T. and Alamukii, N.A.  
725 (2014) Sexual dimorphism in ratio of second and fourth digits and its relationship with metabolic  
726 syndrome indices and cardiovascular risk factors. *Journal of research in medical sciences : the official*  
727 *journal of Isfahan University of Medical Sciences*, **19**, 234-239.

728 21 Teatero, M.L. and Netley, C. (2013) A critical review of the research on the extreme male  
729 brain theory and digit ratio (2D:4D). *Journal of autism and developmental disorders*, **43**, 2664-2676.

730 22 Voracek, M. and Loibl, L.M. (2009) Scientometric analysis and bibliography of digit ratio  
731 (2D:4D) research, 1998-2008. *Psychological reports*, **104**, 922-956.

732 23 McIntyre, M.H. (2006) The use of digit ratios as markers for perinatal androgen action.  
733 *Reproductive biology and endocrinology : RB&E*, **4**, 10.

734 24 Paul, S.N., Kato, B.S., Cherkas, L.F., Andrew, T. and Spector, T.D. (2006) Heritability of the  
735 second to fourth digit ratio (2d:4d): A twin study. *Twin research and human genetics : the official*  
736 *journal of the International Society for Twin Studies*, **9**, 215-219.

737 25 Voracek, M. and Dressler, S.G. (2007) Digit ratio (2D:4D) in twins: heritability estimates and  
738 evidence for a masculinized trait expression in women from opposite-sex pairs. *Psychological*  
739 *reports*, **100**, 115-126.

740 26 Gobrogge, K.L., Breedlove, S.M. and Klump, K.L. (2008) Genetic and environmental  
741 influences on 2D:4D finger length ratios: a study of monozygotic and dizygotic male and female  
742 twins. *Archives of sexual behavior*, **37**, 112-118.

743 27 Medland, S.E. and Loehlin, J.C. (2008) Multivariate genetic analyses of the 2D:4D ratio:  
744 examining the effects of hand and measurement technique in data from 757 twin families. *Twin*  
745 *research and human genetics : the official journal of the International Society for Twin Studies*, **11**,  
746 335-341.

747 28 Voracek, M. and Dressler, S.G. (2009) Brief communication: Familial resemblance in digit  
748 ratio (2D:4D). *American journal of physical anthropology*, **140**, 376-380.

749 29 Lawrance-Owen, A.J., Bargary, G., Bosten, J.M., Goodbourn, P.T., Hogg, R.E. and Mollon, J.D.  
750 (2013) Genetic association suggests that SMOC1 mediates between prenatal sex hormones and digit  
751 ratio. *Human genetics*, **132**, 415-421.

752 30 Medland, S.E., Zayats, T., Glaser, B., Nyholt, D.R., Gordon, S.D., Wright, M.J., Montgomery,  
753 G.W., Campbell, M.J., Henders, A.K., Timpson, N.J. *et al.* (2010) A variant in LIN28B is associated with  
754 2D:4D finger-length ratio, a putative retrospective biomarker of prenatal testosterone exposure.  
755 *American journal of human genetics*, **86**, 519-525.

756 31 He, C., Kraft, P., Chen, C., Buring, J.E., Pare, G., Hankinson, S.E., Chanock, S.J., Ridker, P.M.,  
757 Hunter, D.J. and Chasman, D.I. (2009) Genome-wide association studies identify loci associated with  
758 age at menarche and age at natural menopause. *Nat Genet*, **41**, 724-728.

759 32 Lettre, G., Jackson, A.U., Gieger, C., Schumacher, F.R., Berndt, S.I., Sanna, S., Eyheramendy,  
760 S., Voight, B.F., Butler, J.L., Guiducci, C. *et al.* (2008) Identification of ten loci associated with height  
761 highlights new biological pathways in human growth. *Nat Genet*, **40**, 584-591.

762 33 Abouzeid, H., Boisset, G., Favez, T., Youssef, M., Marzouk, I., Shakankiry, N., Bayoumi, N.,  
763 Descombes, P., Agosti, C., Munier, F.L. *et al.* (2011) Mutations in the SPARC-related modular calcium-  
764 binding protein 1 gene, SMOC1, cause waardenburg anophthalmia syndrome. *American journal of*  
765 *human genetics*, **88**, 92-98.

766 34 Okada, I., Hamanoue, H., Terada, K., Tohma, T., Megarbane, A., Chouery, E., Abou-Ghoch, J.,  
767 Jalkh, N., Cogulu, O., Ozkinay, F. *et al.* (2011) SMOC1 is essential for ocular and limb development in  
768 humans and mice. *American journal of human genetics*, **88**, 30-41.

769 35 Rainger, J., van Beusekom, E., Ramsay, J.K., McKie, L., Al-Gazali, L., Pallotta, R., Saponari, A.,  
770 Branney, P., Fisher, M., Morrison, H. *et al.* (2011) Loss of the BMP antagonist, SMOC-1, causes

771 Ophthalmo-acromelic (Waardenburg Anophthalmia) syndrome in humans and mice. *PLoS Genet*, **7**,  
772 e1002114.

773 36 Love, H.D., Booton, S.E., Boone, B.E., Breyer, J.P., Koyama, T., Revelo, M.P., Shappell, S.B.,  
774 Smith, J.R. and Hayward, S.W. (2009) Androgen regulated genes in human prostate xenografts in  
775 mice: relation to BPH and prostate cancer. *PloS one*, **4**, e8384.

776 37 Schaeffer, E.M., Marchionni, L., Huang, Z., Simons, B., Blackman, A., Yu, W., Parmigiani, G.  
777 and Berman, D.M. (2008) Androgen-induced programs for prostate epithelial growth and invasion  
778 arise in embryogenesis and are reactivated in cancer. *Oncogene*, **27**, 7180-7191.

779 38 Coleman, I.M., Kiefer, J.A., Brown, L.G., Pitts, T.E., Nelson, P.S., Brubaker, K.D., Vessella, R.L.  
780 and Corey, E. (2006) Inhibition of androgen-independent prostate cancer by estrogenic compounds  
781 is associated with increased expression of immune-related genes. *Neoplasia*, **8**, 862-878.

782 39 Berenbaum, S.A., Bryk, K.K., Nowak, N., Quigley, C.A. and Moffat, S. (2009) Fingers as a  
783 marker of prenatal androgen exposure. *Endocrinology*, **150**, 5119-5124.

784 40 Manning, J.T., Bundred, P.E., Newton, D.J. and Flanagan, B.F. (2003) The second to fourth  
785 digit ratio and variation in the androgen receptor gene. *Evolution and Human Behavior*, **24**, 399-405.

786 41 Chamberlain, N.L., Driver, E.D. and Miesfeld, R.L. (1994) The length and location of CAG  
787 trinucleotide repeats in the androgen receptor N-terminal domain affect transactivation function.  
788 *Nucleic acids research*, **22**, 3181-3186.

789 42 Xi, H., Li, M., Fan, Y. and Zhao, L. (2014) A comparison of measurement methods and sexual  
790 dimorphism for digit ratio (2D:4D) in Han ethnicity. *Archives of sexual behavior*, **43**, 329-333.

791 43 Yang, J., Ferreira, T., Morris, A.P., Medland, S.E., Madden, P.A., Heath, A.C., Martin, N.G.,  
792 Montgomery, G.W., Weedon, M.N., Loos, R.J. *et al.* (2012) Conditional and joint multiple-SNP  
793 analysis of GWAS summary statistics identifies additional variants influencing complex traits. *Nat*  
794 *Genet*, **44**, 369-375, s361-363.

795 44 Pers, T.H., Karjalainen, J.M., Chan, Y., Westra, H.J., Wood, A.R., Yang, J., Lui, J.C., Vedantam,  
796 S., Gustafsson, S., Esko, T. *et al.* (2015) Biological interpretation of genome-wide association studies  
797 using predicted gene functions. *Nature communications*, **6**, 5890.

798 45 Jin, G., Sun, J., Kim, S.T., Feng, J., Wang, Z., Tao, S., Chen, Z., Purcell, L., Smith, S., Isaacs, W.B.  
799 *et al.* (2012) Genome-wide association study identifies a new locus JMJD1C at 10q21 that may  
800 influence serum androgen levels in men. *Hum Mol Genet*, **21**, 5222-5228.

801 46 Ohlsson, C., Wallaschowski, H., Lunetta, K.L., Stolk, L., Perry, J.R., Koster, A., Petersen, A.K.,  
802 Eriksson, J., Lehtimäki, T., Huhtaniemi, I.T. *et al.* (2011) Genetic determinants of serum testosterone  
803 concentrations in men. *PLoS Genet*, **7**, e1002313.

804 47 Prescott, J., Thompson, D.J., Kraft, P., Chanock, S.J., Audley, T., Brown, J., Leyland, J., Folkard,  
805 E., Doody, D., Hankinson, S.E. *et al.* (2012) Genome-wide association study of circulating estradiol,  
806 testosterone, and sex hormone-binding globulin in postmenopausal women. *PloS one*, **7**, e37815.

807 48 Coviello, A.D., Haring, R., Wellons, M., Vaidya, D., Lehtimäki, T., Keildson, S., Lunetta, K.L.,  
808 He, C., Fornage, M., Lagou, V. *et al.* (2012) A genome-wide association meta-analysis of circulating  
809 sex hormone-binding globulin reveals multiple Loci implicated in sex steroid hormone regulation.  
810 *PLoS Genet*, **8**, e1002805.

811 49 Bulik-Sullivan, B., Finucane, H.K., Anttila, V., Gusev, A., Day, F.R., Loh, P.R., Duncan, L., Perry,  
812 J.R., Patterson, N., Robinson, E.B. *et al.* (2015) An atlas of genetic correlations across human diseases  
813 and traits. *Nat Genet*, **47**, 1236-1241.

814 50 Kim, Y.S., Lewandoski, M., Perantoni, A.O., Kurebayashi, S., Nakanishi, G. and Jetten, A.M.  
815 (2002) Identification of Glis1, a novel Gli-related, Kruppel-like zinc finger protein containing  
816 transactivation and repressor functions. *The Journal of biological chemistry*, **277**, 30901-30913.

817 51 Eeles, R.A., Olama, A.A., Benlloch, S., Saunders, E.J., Leongamornlert, D.A., Tymrakiewicz, M.,  
818 Ghousaini, M., Luccarini, C., Dennis, J., Jugurnauth-Little, S. *et al.* (2013) Identification of 23 new  
819 prostate cancer susceptibility loci using the iCOGS custom genotyping array. *Nature genetics*, **45**,  
820 385-391, 391e381-382.



821 52 Chambers, J.C., Zhang, W., Sehmi, J., Li, X., Wass, M.N., Van der Harst, P., Holm, H., Sanna, S.,  
822 Kavousi, M., Baumeister, S.E. *et al.* (2011) Genome-wide association study identifies loci influencing  
823 concentrations of liver enzymes in plasma. *Nat Genet*, **43**, 1131-1138.

824 53 Willer, C.J., Schmidt, E.M., Sengupta, S., Peloso, G.M., Gustafsson, S., Kanoni, S., Ganna, A.,  
825 Chen, J., Buchkovich, M.L., Mora, S. *et al.* (2013) Discovery and refinement of loci associated with  
826 lipid levels. *Nat Genet*, **45**, 1274-1283.

827 54 Innocenti, F., Cooper, G.M., Stanaway, I.B., Gamazon, E.R., Smith, J.D., Mirkov, S., Ramirez,  
828 J., Liu, W., Lin, Y.S., Moloney, C. *et al.* (2011) Identification, replication, and functional fine-mapping  
829 of expression quantitative trait loci in primary human liver tissue. *PLoS Genet*, **7**, e1002078.

830 55 Takata, R., Akamatsu, S., Kubo, M., Takahashi, A., Hosono, N., Kawaguchi, T., Tsunoda, T.,  
831 Inazawa, J., Kamatani, N., Ogawa, O. *et al.* (2010) Genome-wide association study identifies five new  
832 susceptibility loci for prostate cancer in the Japanese population. *Nat Genet*, **42**, 751-754.

833 56 Goo, Y.H., Son, S.H., Kreienberg, P.B. and Paul, A. (2014) Novel lipid droplet-associated  
834 serine hydrolase regulates macrophage cholesterol mobilization. *Arteriosclerosis, thrombosis, and*  
835 *vascular biology*, **34**, 386-396.

836 57 Ding, Z., Liu, Y., Rubio, V., He, J., Minze, L.J. and Shi, Z.Z. (2016) OLA1, a Translational  
837 Regulator of p21, Maintains Optimal Cell Proliferation Necessary for Developmental Progression.  
838 *Molecular and cellular biology*, **36**, 2568-2582.

839 58 Adhikari, K., Fuentes-Guajardo, M., Quinto-Sanchez, M., Mendoza-Revilla, J., Camilo Chacon-  
840 Duque, J., Acuna-Alonzo, V., Jaramillo, C., Arias, W., Lozano, R.B., Perez, G.M. *et al.* (2016) A genome-  
841 wide association scan implicates DCHS2, RUNX2, GLI3, PAX1 and EDAR in human facial variation.  
842 *Nature communications*, **7**, 11616.

843 59 Al-Qattan, M.M., Shamseldin, H.E., Salih, M.A. and Alkuraya, F.S. (2017) GLI3-related  
844 polydactyly: a review. *Clinical genetics*, in press.

845 60 Sheth, R., Marcon, L., Bastida, M.F., Junco, M., Quintana, L., Dahn, R., Kmita, M., Sharpe, J.  
846 and Ros, M.A. (2012) Hox genes regulate digit patterning by controlling the wavelength of a Turing-  
847 type mechanism. *Science (New York, N.Y.)*, **338**, 1476-1480.

848 61 Lango Allen, H., Estrada, K., Lettre, G., Berndt, S.I., Weedon, M.N., Rivadeneira, F., Willer,  
849 C.J., Jackson, A.U., Vedantam, S., Raychaudhuri, S. *et al.* (2010) Hundreds of variants clustered in  
850 genomic loci and biological pathways affect human height. *Nature*, **467**, 832-838.

851 62 Okada, Y., Wu, D., Trynka, G., Raj, T., Terao, C., Ikari, K., Kochi, Y., Ohmura, K., Suzuki, A.,  
852 Yoshida, S. *et al.* (2014) Genetics of rheumatoid arthritis contributes to biology and drug discovery.  
853 *Nature*, **506**, 376-381.

854 63 Dubois, P.C., Trynka, G., Franke, L., Hunt, K.A., Romanos, J., Curtotti, A., Zhernakova, A.,  
855 Heap, G.A., Adany, R., Aromaa, A. *et al.* (2010) Multiple common variants for celiac disease  
856 influencing immune gene expression. *Nat Genet*, **42**, 295-302.

857 64 Yang, W., Tang, H., Zhang, Y., Tang, X., Zhang, J., Sun, L., Yang, J., Cui, Y., Zhang, L., Hirankarn,  
858 N. *et al.* (2013) Meta-analysis followed by replication identifies loci in or near CDKN1B, TET3, CD80,  
859 DRAM1, and ARID5B as associated with systemic lupus erythematosus in Asians. *American journal of*  
860 *human genetics*, **92**, 41-51.

861 65 Yang, W., Shen, N., Ye, D.Q., Liu, Q., Zhang, Y., Qian, X.X., Hirankarn, N., Ying, D., Pan, H.F.,  
862 Mok, C.C. *et al.* (2010) Genome-wide association study in Asian populations identifies variants in  
863 ETS1 and WDFY4 associated with systemic lupus erythematosus. *PLoS Genet*, **6**, e1000841.

864 66 Han, J.W., Zheng, H.F., Cui, Y., Sun, L.D., Ye, D.Q., Hu, Z., Xu, J.H., Cai, Z.M., Huang, W., Zhao,  
865 G.P. *et al.* (2009) Genome-wide association study in a Chinese Han population identifies nine new  
866 susceptibility loci for systemic lupus erythematosus. *Nat Genet*, **41**, 1234-1237.

867 67 Dittmer, J. (2003) The biology of the Ets1 proto-oncogene. *Molecular cancer*, **2**, 29.

868 68 Lettice, L.A., Williamson, I., Wiltshire, J.H., Peluso, S., Devenney, P.S., Hill, A.E., Essafi, A.,  
869 Hagman, J., Mort, R., Grimes, G. *et al.* (2012) Opposing functions of the ETS factor family define Shh  
870 spatial expression in limb buds and underlie polydactyly. *Developmental cell*, **22**, 459-467.

69 Kohlhase, J. (2000) SALL1 mutations in Townes-Brocks syndrome and related disorders. *Human mutation*, **16**, 460-466.

70 Townes, P.L. and Brocks, E.R. (1972) Hereditary syndrome of imperforate anus with hand, foot, and ear anomalies. *The Journal of pediatrics*, **81**, 321-326.

71 Kohlhase, J., Hausmann, S., Stojmenovic, G., Dixkens, C., Bink, K., Schulz-Schaeffer, W., Altmann, M. and Engel, W. (1999) SALL3, a new member of the human spalt-like gene family, maps to 18q23. *Genomics*, **62**, 216-222.

72 Wilson, M.G., Towner, J.W., Forsman, I. and Siris, E. (1979) Syndromes associated with deletion of the long arm of chromosome 18[del(18q)]. *American journal of medical genetics*, **3**, 155-174.

73 Zhao, M., Gupta, V., Raj, L., Roussel, M. and Bei, M. (2013) A network of transcription factors operates during early tooth morphogenesis. *Molecular and cellular biology*, **33**, 3099-3112.

74 Kondo, T., Zakany, J., Innis, J.W. and Duboule, D. (1997) Of fingers, toes and penises. *Nature*, **390**, 29.

75 Manning, J.T., Wood, S., Vang, E., Walton, J., Bundred, P.E., van Heyningen, C. and Lewis-Jones, D.I. (2004) Second to fourth digit ratio (2D:4D) and testosterone in men. *Asian journal of andrology*, **6**, 211-215.

76 Lutchmaya, S., Baron-Cohen, S., Raggatt, P., Knickmeyer, R. and Manning, J.T. (2004) 2nd to 4th digit ratios, fetal testosterone and estradiol. *Early Hum Dev*, **77**, 23-28.

77 Hickey, M., Doherty, D.A., Hart, R., Norman, R.J., Mattes, E., Atkinson, H.C. and Sloboda, D.M. (2010) Maternal and umbilical cord androgen concentrations do not predict digit ratio (2D:4D) in girls: a prospective cohort study. *Psychoneuroendocrinology*, **35**, 1235-1244.

78 Voracek, M. (2014) No effects of androgen receptor gene CAG and GGC repeat polymorphisms on digit ratio (2D:4D): a comprehensive meta-analysis and critical evaluation of research. *Evolution and Human Behavior*, **35**, 430-437.

79 Honekopp, J. (2013) No Evidence that 2D:4D is Related to the Number of CAG Repeats in the Androgen Receptor Gene. *Frontiers in endocrinology*, **4**, 185.

80 Wood, A.R., Esko, T., Yang, J., Vedantam, S., Pers, T.H., Gustafsson, S., Chu, A.Y., Estrada, K., Luan, J., Kutalik, Z. *et al.* (2014) Defining the role of common variation in the genomic and biological architecture of adult human height. *Nat Genet*, **46**, 1173-1186.

81 Locke, A.E., Kahali, B., Berndt, S.I., Justice, A.E., Pers, T.H., Day, F.R., Powell, C., Vedantam, S., Buchkovich, M.L., Yang, J. *et al.* (2015) Genetic studies of body mass index yield new insights for obesity biology. *Nature*, **518**, 197-206.

82 Yang, J., Lee, S.H., Goddard, M.E. and Visscher, P.M. (2011) GCTA: a tool for genome-wide complex trait analysis. *American journal of human genetics*, **88**, 76-82.

83 Lee, S.H., Yang, J., Goddard, M.E., Visscher, P.M. and Wray, N.R. (2012) Estimation of pleiotropy between complex diseases using single-nucleotide polymorphism-derived genomic relationships and restricted maximum likelihood. *Bioinformatics (Oxford, England)*, **28**, 2540-2542.

84 Kemp, J.P., Medina-Gomez, C., Estrada, K., St Pourcain, B., Heppe, D.H., Warrington, N.M., Oei, L., Ring, S.M., Kruithof, C.J., Timpson, N.J. *et al.* (2014) Phenotypic dissection of bone mineral density reveals skeletal site specificity and facilitates the identification of novel Loci in the genetic regulation of bone mass attainment. *PLoS Genet*, **10**, e1004423.

85 Medina-Gomez, C., Felix, J.F., Estrada, K., Peters, M.J., Herrera, L., Kruithof, C.J., Duijts, L., Hofman, A., van Duijn, C.M., Uitterlinden, A.G. *et al.* (2015) Challenges in conducting genome-wide association studies in highly admixed multi-ethnic populations: the Generation R Study. *European journal of epidemiology*, **30**, 317-330.

86 Medina-Gomez, C., Kemp, J.P., Dimou, N.L., Kreiner, E., Chesi, A., Zemel, B.S., Bonnelykke, K., Boer, C.G., Ahluwalia, T.S., Bisgaard, H. *et al.* (2017) Bivariate genome-wide association meta-analysis of pediatric musculoskeletal traits reveals pleiotropic effects at the SREBF1/TOM1L2 locus. *Nature communications*, **8**, 121.

87 Kappen, J.H., Medina-Gomez, C., van Hagen, P.M., Stolk, L., Estrada, K., Rivadeneira, F.,  
 88 Uitterlinden, A.G., Stanford, M.R., Ben-Chetrit, E., Wallace, G.R. *et al.* (2015) Genome-wide  
 89 association study in an admixed case series reveals IL12A as a new candidate in Behcet disease. *PloS*  
 90 *one*, **10**, e0119085.  
 91 van der Valk, R.J., Duijts, L., Timpson, N.J., Salam, M.T., Standl, M., Curtin, J.A., Genuneit, J.,  
 92 Kerhof, M., Kreiner-Moller, E., Caceres, A. *et al.* (2014) Fraction of exhaled nitric oxide values in  
 93 childhood are associated with 17q11.2-q12 and 17q12-q21 variants. *The Journal of allergy and*  
 94 *clinical immunology*, **134**, 46-55.  
 95 Willer, C.J., Li, Y. and Abecasis, G.R. (2010) METAL: fast and efficient meta-analysis of  
 96 genomewide association scans. *Bioinformatics (Oxford, England)*, **26**, 2190-2191.  
 97 Allen, N.E., Sudlow, C., Peakman, T. and Collins, R. (2014) UK biobank data: come and get it.  
 98 *Science translational medicine*, **6**, 224ed224.  
 99 GETEx Consortium. (2013) The Genotype-Tissue Expression (GTEx) project. *Nat Genet*, **45**,  
 100 580-585.  
 101 Boker, S., Neale, M., Maes, H., Wilde, M., Spiegel, M., Brick, T., Spies, J., Estabrook, R.,  
 102 Kenny, S., Bates, T. *et al.* (2011) OpenMx: An Open Source Extended Structural Equation Modeling  
 103 Framework. *Psychometrika*, **76**, 306-317.  
 104 Ihaka, R., Gentleman R. (1996) R: a language for data analysis and graphics. *Journal of*  
 105 *Computational and Graphical Statistics*, **5**, 299-314.

**Figure Legends:**

**Figure 1:** Manhattan plot from the discovery meta-analysis of left hand 2D:4D ratio. The red line indicates genome-wide significance ( $P < 5 \times 10^{-8}$ ) and the yellow line indicates suggestive significance ( $P < 1 \times 10^{-5}$ ). Purple dots indicate those loci that reach genome-wide significance.

**Figure 2:** Plots highlighting the relationship of 13 SHBG and five testosterone associated SNPs (45-48) with the left hand 2D:4D ratio discovery meta-analysis. (A) Q-Q plot of the meta-analysis P-values for each of the SNPs. (B) Plot of the  $\beta$  coefficient for the left hand 2D:4D ratio meta-analysis against the  $\beta$  coefficient previously reported for SHBG or testosterone (the 13 SHBG associated SNPs are aligned to the SHBG increasing allele and the five testosterone associated SNPs are aligned to the testosterone increasing allele). Triangles indicate SNPs associated with testosterone, circles indicate SNPs associated with SHBG.

955 **Table 1:** Descriptive statistics of the discovery and replication cohorts.

Variable	Subset	ALSPAC	Generation R	QIMR	Raine	Rotterdam Study	Twins UK
N	All	5,337	3,059	2,775	1,003	2,091	1,396
Age (years) <sup>a</sup>	All	11.74 (0.23)	9.80 (0.33)	15.47 (2.93)	20.05 (0.43)	67.84 (7.91)	54.84 (12.21)
Sex (male) <sup>b</sup>	All	49% (2615)	47.9% (1465)	46.4% (1287)	50.85% (510)	42.9% (897)	9.10% (127)
Left 2D:4D <sup>a</sup>	All	96.53 (3.25)	91.15 (2.72)	97.66 (3.41)	96.52 (3.45)	92.40 (2.24)	96.70 (3.40)
	Male	96.05 (3.17)	90.74 (2.77)	96.87 (3.38)	96.16 (3.33)	91.96 (2.25)	95.25 (3.24)
	Female	97.00 (3.26)	91.52 (2.62)	98.33 (3.29)	96.93 (3.53)	92.73 (2.17)	96.80 (3.42)
Right 2D:4D <sup>a</sup>	All	96.37 (3.28)	--	97.07 (3.43)	96.99 (3.28)	92.39 (2.42)	97.10 (3.50)
	Male	95.87 (3.22)	--	96.12 (3.29)	96.81 (3.28)	91.94 (2.39)	95.58 (3.38)
	Female	96.87 (3.26)	--	97.88 (3.33)	97.19 (3.28)	92.73 (2.39)	97.24 (3.48)
Mean 2D:4D <sup>a</sup>	All	96.45 (2.99)	--	97.39 (3.10)	96.76 (3.04)	92.40 (2.12)	96.90 (3.10)
	Male	95.96 (2.90)	--	96.51 (3.01)	96.48 (2.94)	91.95 (2.13)	95.42 (2.99)
	Female	96.93 (3.00)	--	98.14 (2.97)	97.06 (3.12)	92.71 (2.05)	97.02 (3.07)

956 <sup>a</sup> Mean (SD); <sup>b</sup> Percent (number)

**Table 2:** Genome-wide-significant loci from the discovery meta-analysis in all individuals for left hand 2D:4D ratio; the most significant SNP from each locus is presented. Replication results are presented from 23andMe where the 2D:4D ratio was reported as a relative measure (i.e. 0 = index finger longer [17.1% of research participants], 1 = index and ring finger the same length [14.0% of research participants], 2 = ring finger longer [68.9% of research participants]).

	Chr	Position (bp [GRCh37/ hg19])	Nearest gene	Effect allele / Other allele	EAFA	Beta	SE	P-Value
rs4927012								
Discovery	1	54068016	<i>GLIS1</i>	T/C	0.875	-0.358	0.058	5.08x10 <sup>-10</sup>
Replication	1	54068016	<i>GLIS1</i>	T/C	0.871	-0.042	0.006	3.48x10 <sup>-12</sup>
rs11581730								
Discovery	1	155082158	<i>EFNA1</i>	A/T	0.496	0.294	0.036	3.02x10 <sup>-16</sup>
Replication	1	155082158	<i>EFNA1</i>	A/T	0.503	0.026	0.004	4.99x10 <sup>-11</sup>
rs340600								
Discovery	2	20892006	<i>LDAH<sup>b</sup></i>	T/G	0.199	-0.379	0.046	1.38x10 <sup>-16</sup>
Replication	2	20892006	<i>LDAH<sup>b</sup></i>	T/G	0.199	-0.043	0.005	1.81x10 <sup>-17</sup>
rs12474669								
Discovery	2	175134232	<i>OLA1</i>	A/G	0.139	0.417	0.054	1.51x10 <sup>-14</sup>
Replication	2	175134232	<i>OLA1</i>	A/G	0.143	0.043	0.006	1.92x10 <sup>-13</sup>
rs847158								
Discovery	2	176962102	<i>HOXD12/HOXD11</i>	A/G	0.602	0.199	0.037	1.03x10 <sup>-7</sup>
Replication	2	176962102	<i>HOXD12/HOXD11</i>	A/G	0.609	0.039	0.004	3.04x10 <sup>-20</sup>
rs314277 <sup>c</sup>								
Discovery	6	105407662	<i>LIN28B</i>	A/C	0.155	0.428	0.050	5.55x10 <sup>-18</sup>
Replication	6	105407662	<i>LIN28B</i>	A/C	0.149	0.067	0.006	1.77x10 <sup>-32</sup>
rs77640775 <sup>d</sup>								
Discovery	7	42190714	<i>GLI3</i>	A/G	0.137	-0.252	0.053	1.92x10 <sup>-6</sup>
Replication	7	42190714	<i>GLI3</i>	A/G	0.146	-0.033	0.006	6.64x10 <sup>-9</sup>
rs10790969								
Discovery	11	128529842	<i>FLI1</i>	T/C	0.276	0.284	0.040	1.33x10 <sup>-12</sup>
Replication	11	128529842	<i>FLI1</i>	T/C	0.272	0.027	0.005	1.26x10 <sup>-9</sup>
rs2332175 <sup>c</sup>								
Discovery	14	70345411	<i>SMOC1</i>	A/G	0.529	0.360	0.037	3.00x10 <sup>-22</sup>
Replication	14	70345411	<i>SMOC1</i>	A/G	0.546	0.045	0.004	6.74x10 <sup>-29</sup>
rs6499762								
Discovery	16	51697874	<i>SALL1</i>	A/C	0.125	0.441	0.056	5.33x10 <sup>-15</sup>
Replication	16	51697874	<i>SALL1</i>	A/C	0.129	0.083	0.006	2.83x10 <sup>-41</sup>
rs1080014								
Discovery	16	51900171	<i>TOX3</i>	C/T	0.514	0.203	0.036	1.94x10 <sup>-8</sup>

	Chr	Position (bp [GRCh37/ hg19])	Nearest gene	Effect allele / Other allele	EAF <sup>a</sup>	Beta	SE	P-Value
Replication	16	51900171	<i>TOX3</i>	C/T	0.501	0.012	0.004	3.35x10 <sup>-3</sup>
rs4799176								
Discovery	18	76378307	<i>SALL3</i>	C/T	0.256	0.305	0.044	4.09x10 <sup>-12</sup>
Replication	18	76378307	<i>SALL3</i>	C/T	0.244	0.057	0.005	1.89x10 <sup>-34</sup>

<sup>a</sup> Average effect allele frequency (EAF) across the cohorts in each of the meta-analyses.

<sup>b</sup> Previously known as *C2orf43*.

<sup>c</sup> Genetic loci that had previously been associated with 2D:4D ratio in Medland et al. 2010.

<sup>d</sup> SNP passed genome-wide significance in the average 2D:4D ratio meta-analysis (see supplementary material for results)

**Table 3:** Association between the number of CAG repeats in the *AR* gene and the mean of the left and right hand 2D:4D ratios. Displayed are beta (SE) and P-values in each of the cohorts and the combined estimates from the fixed effects, inverse-variance weighted meta-analysis. ‘Mean’ refers to analyses involving the average CAG repeat length, ‘High’ refers to analyses involving the highest length repeat and ‘Low’ refers to analyses involving the lower length repeat. One-tailed P-values testing for a positive association between CAG repeat length and 2D:4D ratio are presented.

	ALSPAC	QIMR	Meta-Analysis
<i>All Individuals</i>			
	N=5328	N=498	N=5826
Mean	0.014 (0.016), P=0.19	0.040 (0.052), P=0.22	0.016 (0.015), P=0.14
High	0.012 (0.015), P=0.21	0.014 (0.046), P=0.38	0.012 (0.014), P=0.20
Low	0.013 (0.016), P=0.21	0.041 (0.057), P=0.24	0.015 (0.015), P=0.16
<i>Male</i>			
	N=2615	N=231	N=2846
Mean	-0.002 (0.020), P=0.54	-0.099 (0.083), P=0.88	-0.007 (0.019), P=0.65
<i>Female</i>			
	N=2713	N=287	N=3000
Mean	0.046 (0.028), P=0.05	0.125 (0.072), P=0.04	0.056 (0.026), P=0.02
High	0.030 (0.025), P=0.12	0.068 (0.058), P=0.12	0.036 (0.023), P=0.06
Low	0.040 (0.026), P=0.06	0.112 (0.082), P=0.09	0.047 (0.025), P=0.03



**Supplementary Information Captions:**

Supplementary information includes:

**S1-Figure:** QQ Plots from the discovery GWAS meta-analysis for left hand (European only; A), left hand (Multiethnic; B), right hand (C) and mean (D) 2D:4D ratio.

**S2-Figure:** QQ Plots of the heterogeneity P-value from the discovery meta-analysis for the left hand (European only; A), left hand (Multi-ethnic; B), right hand (C) and mean (D) 2D:4D ratio.

**S3-Figure:** Manhattan plots from the discovery meta-analysis for the left hand (multiethnic; A), right hand (B) and average of both hands (C) 2D:4D ratios. The red line indicates genome-wide significance ( $P < 5 \times 10^{-8}$ ) and the yellow line indicates suggestive significance ( $P < 1 \times 10^{-5}$ ). Purple dots indicate those loci that reach genome-wide significance.

**S4-Figure:** Regional association plot for each of the loci reaching genome-wide significance with the top SNP highlighted in purple.

**S5-Figure:** Region plots from the discovery meta-analysis of the *HOXA* gene cluster and 200kb either side of the gene cluster, for the left hand (European; A), left hand (Multiethnic; B) right hand (C) and mean (D) 2D:4D ratio. The SNP with the lowest P-value in the region for each of the three phenotypes is highlighted in purple.

**S6-Figure:** Miami plots from the sex-stratified analyses, with females on the upper axis and males on the lower, for the left hand (European; A), left hand (Multiethnic; B), right hand (C) and average of both hands (D) 2D:4D ratios. The red line indicates genome-wide significance ( $P < 5 \times 10^{-8}$ ) and the yellow

line indicates suggestive significance ( $P < 1 \times 10^{-5}$ ). Purple dots indicate those loci that reach genome-wide significance.

**S7-Figure:** Manhattan plots from the heterogeneity test between male and female effect sizes for the left hand (European; A), left hand (Multiethnic; B), right hand (C) and average of both hands (D) 2D:4D ratios. The red line indicates genome-wide significance ( $P < 5 \times 10^{-8}$ ) and the yellow line indicates suggestive significance ( $P < 1 \times 10^{-5}$ ). Purple dots indicate those loci that reach genome-wide significance.

**S8-Figure:** QQ Plots from the heterogeneity test between male and female effect sizes for left hand (European only; A), left hand (Multiethnic; B), right hand (C) and mean (D) 2D:4D ratio.

**S9-Figure:** Gene expression profiles of the nearest gene to the lead SNPs using data from the GTEx Consortium. The y-axis shows the Reads Per Kilobase of transcript per Million mapped reads (RPKM).

**S10-Figure:** Region plots from the discovery meta-analysis of the *AR* gene, and 200kb either side of the gene, for the left hand (European; A), left hand (Multiethnic; B) right hand (C) and mean (D) 2D:4D ratio. The SNP with the lowest P-value in the region for each of the three phenotypes is highlighted in purple.

**S1-Table:** Reliability of the CAG(n) polymorphism in the ALSPAC cohort across genotyping replicates in the subset of individuals used for quality control (N=370).

**S2-Table:** Power, bias and coverage probability results from the simulations mimicking the ALSPAC data where the number of CAG repeats was simulated without and with measurement error.

**S3-Table:** Genome-wide-significant loci from the discovery meta-analysis in all individuals; the most significant SNP from each locus is presented.

**S4-Table:** Most likely causal gene at each locus identified by DEPICT

**S5-Table:** Results from the geneset enrichment analysis in DEPICT

**S6-Table:** Results from the tissue enrichment analysis in DEPICT

**S7-Table:** Results for the regression of 2D:4D ratio (left and right hand) on number of CAG repeats in the *AR* gene. Each cell displays the beta coefficient (SE) and P-value from the regression. Results are presented for each of the cohorts and the combined estimates from the fixed effects, inverse-variance weighted meta-analysis. 'Mean' is the average repeat length, 'High' the highest length repeat and 'Low' is the lower length repeat.

**S8-Table:** Genome-wide SNP-heritability and genetic correlation between 2D:4D ratio and a range of traits and diseases.

**Author Contributions:**

N.M.W., E.S., P.G.H., Y.J., C.M-G., K.T., S.E.M. and D.M.E. performed study-level data analysis and N.M.W. performed all meta-analyses. G.H., P.G.H., C.A.W., J.P.K., G.M., D.W., D.A.M., N.G.M., G.D.S., C.E.P. and T.D.S. performed study-level genotyping and imputation. Replication services provided by Y.J. and A.A., and the 23andMe Research Team. Study design was by S.E.M. and D.M.E. Data collection and interpretation was by C.G.B., M.H., M.A.I., the 23andMe Research team, G.W.M., J.F.F., M.J.W., D.A.M., V.W.J., N.G.M, J.Y.T., G.D.S., C.E.P., T.D.S., J.V.M., F.R., S.E.M. and D.M.E. The paper was written by N.M.W. and D.M.E., and all authors reviewed and revised the paper