

1 **Phenotypic clustering reveals distinct subtypes of polycystic ovary**

2 **syndrome with novel genetic associations**

3

4 Matthew Dapas¹, Frederick T. J. Lin¹, Girish N. Nadkarni², Ryan Sisk¹, Richard S. Legro³,

5 Margrit Urbanek^{1,4,5}, M. Geoffrey Hayes^{1,4,6¶*}, Andrea Dunaif^{7,¶*}

6

7 ¹ Division of Endocrinology, Metabolism, and Molecular Medicine, Department of Medicine,

8 Northwestern University Feinberg School of Medicine, Chicago, IL

9 ² Division of Nephrology, Icahn School of Medicine at Mount Sinai, New York, NY

10 ³ Department of Obstetrics and Gynecology, Penn State College of Medicine, Hershey, PA

11 ⁴ Center for Genetic Medicine, Northwestern University Feinberg School of Medicine, Chicago, IL

12 ⁵ Center for Reproductive Science, Northwestern University Feinberg School of Medicine, Chicago, IL

13 ⁶ Department of Anthropology, Northwestern University, Evanston, IL

14 ⁷ Division of Endocrinology, Diabetes and Bone Disease, Icahn School of Medicine at Mount Sinai, New
15 York, NY

16

17 ¶ These authors jointly supervised this work.

18 * Corresponding author

19

20 E-mail: andrea.dunaif@mssm.edu (AD)

21 ghayes@northwestern.edu (MGH)

22

23 DISCLOSURE STATEMENT: The authors have nothing to disclose.

24 **Abstract**

25 **Background**

26 Polycystic ovary syndrome (PCOS) is a common, complex genetic disorder affecting up to 15%
27 of reproductive age women worldwide, depending on the diagnostic criteria applied. These
28 diagnostic criteria are based on expert opinion and have been the subject of considerable
29 controversy. The phenotypic variation observed in PCOS is suggestive of an underlying genetic
30 heterogeneity, but a recent meta-analysis of European ancestry PCOS cases found that the
31 genetic architecture of PCOS defined by different diagnostic criteria was generally similar,
32 suggesting that the criteria do not identify biologically distinct disease subtypes. We performed
33 this study to test the hypothesis that there are biologically relevant subtypes of PCOS.

35 **Methods and Findings**

36 Unsupervised hierarchical cluster analysis was performed on quantitative anthropometric,
37 reproductive, and metabolic traits in a genotyped discovery cohort of 893 PCOS cases and an
38 ungenotyped validation cohort of 263 PCOS cases. We identified two PCOS subtypes: a
39 “reproductive” group (21-23%) characterized by higher luteinizing hormone (LH) and sex
40 hormone binding globulin (SHBG) levels with relatively low body mass index (BMI) and
41 insulin levels; and a “metabolic” group (37-39%), characterized by higher BMI, glucose, and
42 insulin levels with lower SHBG and LH levels. We performed a GWAS on the genotyped
43 cohort, limiting the cases to either the reproductive or metabolic subtypes. We identified
44 alleles in four novel loci that were associated with the reproductive subtype at genome-wide
45 significance (*PRDM2/KAZNI*, $P=2.2\times 10^{-10}$; *IQCA1*, $P=2.8\times 10^{-9}$; *BMPR1B/UNC5C*,
46 $P=9.7\times 10^{-9}$; *CDH10*, $P=1.2\times 10^{-8}$) and one locus that was significantly associated with the

47 metabolic subtype (*KCNH7/FIGN*, $P=1.0\times 10^{-8}$). We have previously reported that rare
48 variants in *DENNDIA*, a gene regulating androgen biosynthesis, were associated with PCOS
49 quantitative traits in a family-based whole genome sequencing analysis. We classified the
50 reproductive and metabolic subtypes in this family-based PCOS cohort and found that the
51 subtypes tended to cluster in families and that carriers of rare *DENNDIA* variants were
52 significantly more likely to have the reproductive subtype of PCOS. Limitations of our study
53 were that only PCOS cases of European ancestry diagnosed by NIH criteria were included,
54 the sample sizes for the subtype GWAS were small, and the GWAS findings were not
55 replicated.

56

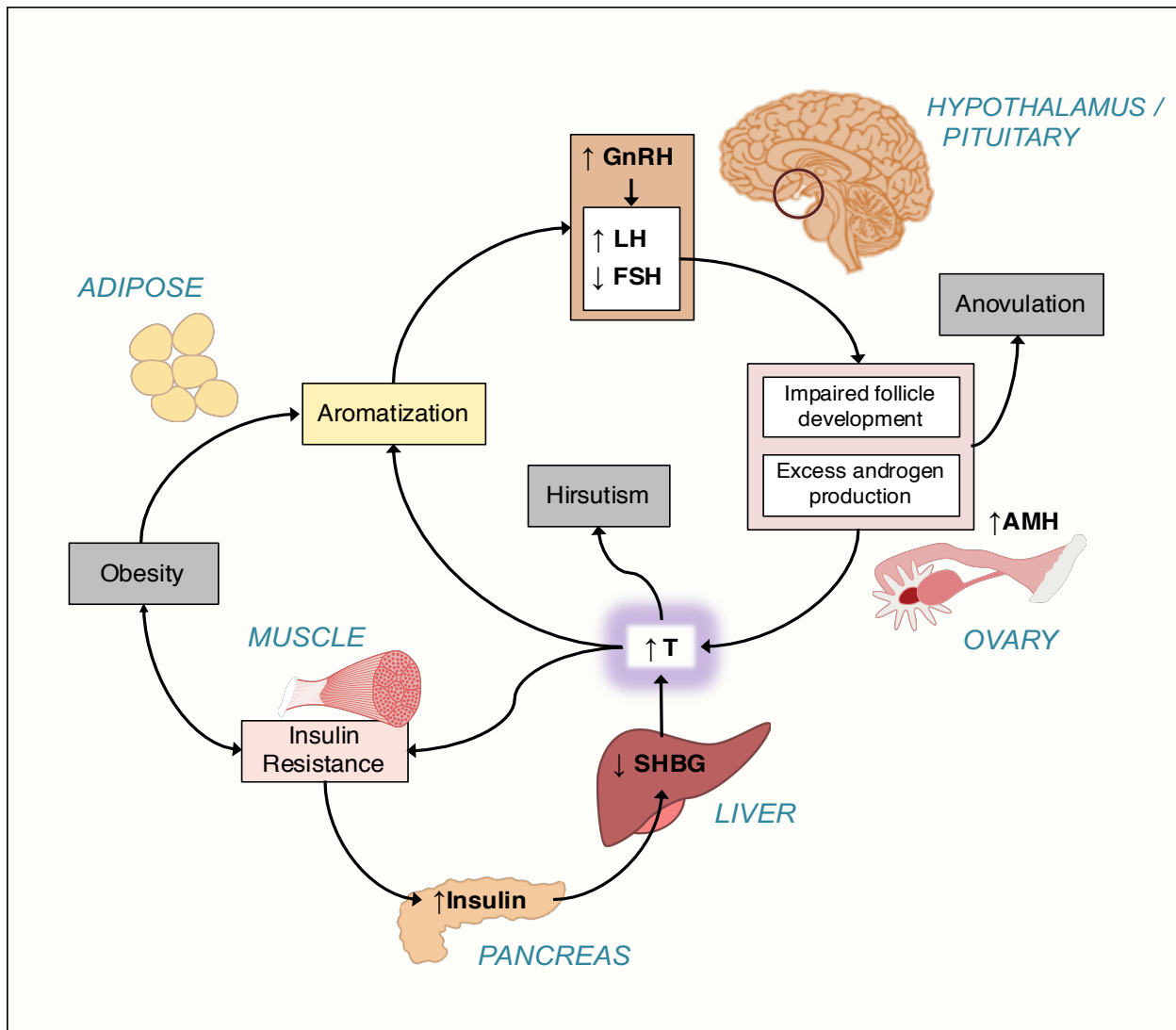
57 **Conclusions**

58 In conclusion, we have found stable reproductive and metabolic subtypes of PCOS. Further,
59 these subtypes were associated with novel susceptibility loci. Our results suggest that these
60 subtypes are biologically relevant since they have distinct genetic architectures. This study
61 demonstrates how precise phenotypic delineation can be more powerful than increases in
62 sample size for genetic association studies.

63 **Introduction**

64 Understanding the genetic architecture of complex diseases is a central challenge in
65 human genetics (1-3). Often defined according to arbitrary diagnostic criteria, complex
66 diseases can represent the phenotypic convergence of numerous genetic etiologies (4-8).
67 Recent studies in type 2 diabetes support the concept that there are disease subtypes with
68 distinct genetic architecture (7,8). Identifying and addressing genetic heterogeneity in
69 complex diseases could increase power to detect causal variants and improve treatment
70 efficacy (9).

71 Polycystic ovary syndrome (PCOS) is a highly heritable, complex genetic disorder
72 affecting up to 15% of reproductive-age women worldwide, depending on the diagnostic
73 criteria applied (10). It is characterized by a variable constellation of reproductive and
74 metabolic abnormalities (**Fig 1**). It is the leading cause of anovulatory infertility and a major
75 risk factor for type 2 diabetes (T2D) in women (11). Despite these substantial morbidities,
76 the etiology(ies) of PCOS remains unknown (12). Accordingly, the commonly-used
77 diagnostic criteria for PCOS, the National Institutes of Health (NIH) criteria (13) and the
78 Rotterdam criteria (14,15), are based on expert opinion, rather than mechanistic insights, and
79 are designed to account for the diverse phenotypic presentations of PCOS. The NIH criteria
80 require the presence of hyperandrogenism (HA) and chronic oligo/anovulation or ovarian
81 dysfunction (OD) (13). The Rotterdam criteria include polycystic ovarian morphology
82 (PCOM) and require the presence of at least two of these three key reproductive traits,
83 resulting in three different affected phenotypes: HA and OD with or without PCOM, also
84 known as NIH PCOS, as well as two additional non-NIH Rotterdam phenotypes, HA and
85 PCOM, and OD and PCOM.



86

87 **Fig 1. Pathophysiology of PCOS** (reviewed in (10)). There is increased frequency of
 88 pulsatile gonadotropin-releasing hormone (GnRH) secretion from the arcuate nucleus of
 89 the hypothalamus that selectively increases luteinizing hormone (LH) secretion. LH
 90 stimulates ovarian theca cell testosterone (T) production. T is incompletely aromatized
 91 to estradiol by the adjacent granulosa cells because of relative follicle stimulating
 92 hormone (FSH) deficiency. There are constitutive increases in the activity of multiple
 93 steroidogenic enzymes in theca cells from women with PCOS, which contributes to
 94 increased T production. Increased adrenal androgen production may also be present in
 95 PCOS. T acts in the periphery to produce clinical signs of androgen excess, such as
 96 hirsutism, acne, and alopecia. T and androstenedione can also be aromatized
 97 extragonadally to estradiol and estrone, respectively, resulting in unopposed estrogen
 98 action on the endometrium. T feeds back on the hypothalamus to decrease the
 99 sensitivity to the normal feedback effects of estradiol and progesterone to slow GnRH
 100 pulse frequency. Anti-Müllerian hormone (AMH) levels are frequently increased in

101 PCOS; this hormone is secreted by small, growing preantral follicles, which are
102 increased in PCOS. Recent studies suggest AMH acting through its cognate receptor on
103 GnRH neurons in the arcuate nucleus contributes to the pathogenesis of PCOS (16,17).

104 PCOS is often associated with profound insulin resistance due to a unique defect
105 in post-binding insulin-mediated signal transduction. Insulin is a co-gonadotropin that
106 acts in synergy with LH to amplify the reproductive abnormalities of PCOS. Insulin
107 signaling in the hypothalamus also appears to be important for ovulation. Insulin is a
108 major negative regulator of hepatic synthesis of sex hormone-binding globulin (SHBG),
109 the specific transport protein for T; only T which is not bound to SHBG is biologically
110 active.

111

112 Genomewide association studies (GWAS) have considerably advanced our
113 understanding of the pathophysiology of PCOS. These studies have implicated gonadotropin
114 secretion (18) and action (19,20), androgen biosynthesis (19-21), metabolic
115 regulation (21,22) and ovarian aging (22) in PCOS pathogenesis. A recent meta-analysis (21)
116 of GWAS was the first study to investigate the genetic architecture of the diagnostic criteria.
117 Only one of 14 PCOS susceptibility loci identified was significantly more strongly associated
118 with the NIH phenotype compared to non-NIH Rotterdam phenotypes or to self-reported
119 PCOS. These findings suggested that the genetic architecture of the phenotypes defined by
120 the different PCOS diagnostic criteria was generally similar. Therefore, the current
121 diagnostic criteria do not appear to identify genetically distinct disease subtypes.

122 It is possible to identify physiologically relevant complex disease subtypes through
123 cluster analysis of phenotypic traits (7,23,24). Indeed, there have been previous efforts to
124 subtype PCOS using unsupervised cluster analysis of its hormonal and anthropometric
125 traits (25-28). However, there has been no validation that the resulting PCOS subtypes were
126 biologically meaningful by testing their association with genetic variants, with other
127 independent biomarkers, or with outcomes, such as therapeutic responses. In this study, we
128 sought to 1) identify phenotypic subtypes of PCOS using an unsupervised clustering

129 approach on reproductive and metabolic quantitative traits from a large cohort of women
130 with PCOS, 2) validate those subtypes in a replication cohort, and 3) test whether the
131 subtypes thus identified were associated with distinct common genetic variants. As an
132 additional validation, we investigated the association of the subtypes with rare genetic
133 variants we recently identified in a family-based PCOS cohort (29).

134

135 **Methods**

136 **Subjects**

137 This study used biochemical and genotype data from the previously published PCOS
138 GWAS, Hayes and Urbanek et al., 2015 (18), in which a discovery cohort (Stage 1) of 984
139 PCOS cases and 2,964 population controls were studied, followed by a replication cohort
140 (Stage 2) of 1,799 PCOS cases and 1,231 phenotyped reproductively normal control women.
141 All cases were of European ancestry and each subject provided written informed consent
142 prior to the study (18). PCOS cases were ages 13-45 years and were diagnosed according to
143 the NIH criteria (10) of hyperandrogenism and chronic anovulation (eight or fewer menses
144 per year), excluding specific disorders of the adrenals, ovaries, or pituitary (30). Cases
145 fulfilling the NIH criteria also meet the Rotterdam criteria for PCOS (10). Two additional
146 PCOS cohorts were included in the present study who fulfilled the NIH criteria and were
147 phenotyped according to the same methods as the genotyped GWAS cohort. An independent,
148 ungenotyped cohort of 263 women with PCOS was used for clustering replication. A family-
149 based whole-genome sequencing cohort of 73 women with PCOS was used for investigating
150 subtype clustering in families and for rare variant analysis (29).

151 Population-based control DNA samples for the GWAS Stage 1 cohort were obtained
152 from the NUGene biobank (31) from women of European ancestry, ages 18-97 years. Control
153 women in the Stage 2 cohort were phenotyped reproductively normal women of European
154 ancestry, ages 15–45 years, with regular menses and normal T levels, and who were not
155 receiving contraceptive steroids for at least 3 months prior to study (18). T, DHEAS, SHBG,
156 luteinizing hormone (LH), follicle-stimulating hormone (FSH), fasting glucose (Glu0), and
157 fasting insulin (Ins0) levels were measured as previously reported (18).

158

159 **Clustering**

160 Clustering was performed in PCOS cases on eight adjusted quantitative traits: BMI,
161 T, DHEAS, Ins0, Glu0, SHBG, LH, and FSH. There were 893 combined cases from both
162 stages with complete quantitative trait data available for clustering (**Table S1**). Quantitative
163 trait values were first log_e-normalized and adjusted for age and assay method, which varied
164 according to the different study sites where samples were collected (18), using a linear
165 regression. An inverse normal transformation was then applied for each trait to ensure equal
166 scaling. The normalized trait residuals were clustered using unsupervised, agglomerative,
167 hierarchical clustering according to a generalization of Ward's minimum variance
168 method (32,33) on Manhattan distances between trait values. Differences in adjusted,
169 normalized trait values between subtypes were assessed using Kruskal-Wallis and pairwise
170 Wilcoxon rank-sum tests corrected for multiple testing (Bonferroni). Cluster stability was
171 assessed by computing the mean Jaccard coefficient from a repeated nonparametric bootstrap
172 resampling (n=1000) of the dissimilarity matrix (34). Jaccard coefficients below 0.5 indicate
173 that a cluster does not capture any discernable pattern within the data, while a mean

174 coefficient above 0.6 indicates that the cluster reflects a real pattern within the data (34).

175 Cluster reproducibility was further assessed by repeating the clustering procedure in an
176 independent cohort of 263 PCOS cases.

177

178 **Association Testing**

179 Stage 1 samples were genotyped using the Illumina OmniExpress
180 (HumanOmniExpress-12v1_C) single nucleotide polymorphism (SNP) array. Stage 2
181 samples were genotyped using the MetaboChip (35) with added custom variant content based
182 on ancillary studies and the discovery results (18). The Stage 2 association replication in this
183 study was therefore limited; many of the loci from Stage 1 were therefore not characterized in
184 Stage 2. Low quality genotypic data were removed as described previously (18). SNPs were
185 filtered according to minor allele frequency ($MAF \geq 0.01$), Hardy-Weinberg equilibrium (p
186 $\geq 1 \times 10^{-6}$), call rate (≥ 0.99), minor allele count ($MAC > 5$), Mendelian concordance, and
187 duplicate sample concordance. Only autosomal SNPs were considered. Ancestry was
188 evaluated using a principal component analysis (PCA) (36) on 76,602 LD-pruned SNPs (18).
189 Samples with values > 3 standard deviations from the median for either of the first two
190 principal components (PCs) were excluded (34 in discovery; 37 in replication). Genotype
191 data was phased using ShapeIT (v2.r790) (37) and then imputed to the 1000 Genomes
192 reference panel (Phase3 v5) (38) using Minimac3 via the Michigan Imputation Server (39).
193 Imputed SNPs with an allelic r^2 below 0.8 were removed from analysis.

194 Association testing was performed separately for Stage 1 and Stage 2 cohorts. Of the
195 893 combined cases from both stages included in the clustering analysis, 555 were from the
196 Stage 1 cohort and 338 were from the Stage 2 cohort. 2,964 normal controls were used in

197 Stage 1, and 1,134 were used in Stage 2. Logistic regressions were performed using
198 SNPTEST (40) for case-control status under an additive genetic model, adjusting for BMI
199 and first three PCs of ancestry. P-values are reported as P_1 and P_2 for Stage 1 and Stage 2,
200 respectively. Cases were limited to specific subtypes selected from clustering results. The
201 betas and standard errors were combined across Stage 1 and Stage 2 cohorts for each subtype
202 under a fixed meta-analysis model weighting each strata by sample size (41). Association test
203 outputs were aligned to the same reference alleles and weighted z-scores were computed for
204 each SNP. The square root of the cohort-specific sample size was used as the proportional
205 weight. Meta-analysis P-values (P_{meta}) were adjusted for genomic inflation. Associations with
206 P-values $< 1.67 \times 10^{-8}$ were considered statistically significant, based on the standard $P < 5 \times 10^{-8}$
207 used in conventional GWAS adjusted for the three independent association tests performed.

208

209 **Chromatin interactions**

210 Neighboring chromatin interactions were investigated in intergenic loci using high-
211 throughput chromatin conformation capture (Hi-C) data from the 3DIV database (42).
212 Topologically associating domains (TADs) were identified using TopDom (43) with a
213 window size of 20.

214

215 **Identifying subtypes in PCOS families**

216 Quantitative trait data from the affected women ($n=73$) in the family-sequencing
217 cohort (29) were adjusted and normalized as described above. Subtype classifiers were
218 modeled on the adjusted trait values and cluster assignments from the genotyped cohort.
219 Several classification methods were compared using 10-fold cross-validation, including

220 support vector machine, random forest (RF), Gaussian mixed-model, and quadratic
221 discriminant analysis (44). The classifier with the lowest error rate was then applied to the
222 affected women in the family-sequencing cohort to identify subtypes of PCOS in the family
223 data. Some of the probands from the family-based cohort were included in our previous
224 GWAS (18). Therefore, there was some sample overlap between the training and test
225 cohorts: of the 893 genotyped women used to identify the original subtype clusters, 47 were
226 also probands in the family-based cohort. Differences between subtypes in the proportion of
227 women with *DENNDIA* rare variants were tested using the chi-square test of independence.

228

229 **Results**

230 **PCOS subtypes**

231 The clustering revealed two distinct phenotypic subtypes: 1) a group (23%)
232 characterized by higher LH and SHBG levels with relatively low BMI and Ins0 levels, which
233 we designated “reproductive”, and 2) a group (37%) characterized by higher BMI and Glu0
234 and Ins0 levels with relatively low SHBG and LH levels, which we designated “metabolic”
235 (**Fig 2**). The key traits distinguishing the reproductive and metabolic subtypes were BMI,
236 insulin, SHBG, glucose, LH, and FSH, in order of importance according to relative pairwise
237 Wilcoxon rank-sum test statistics (**Fig 3**). The remaining cases (40%) demonstrated no
238 distinguishable pattern regarding their relative phenotypic trait distributions and were
239 designated “indeterminate”. The reproductive and metabolic subtypes clustered along
240 opposite ends of the SHBG vs. Ins0/BMI axis, which was highly correlated with the first PC
241 of the adjusted quantitative traits (**Fig 4**). The reproductive subtype was the most stable
242 cluster, with a mean bootstrapped Jaccard coefficient ($\bar{\gamma}_C$) of 0.61, followed by the metabolic

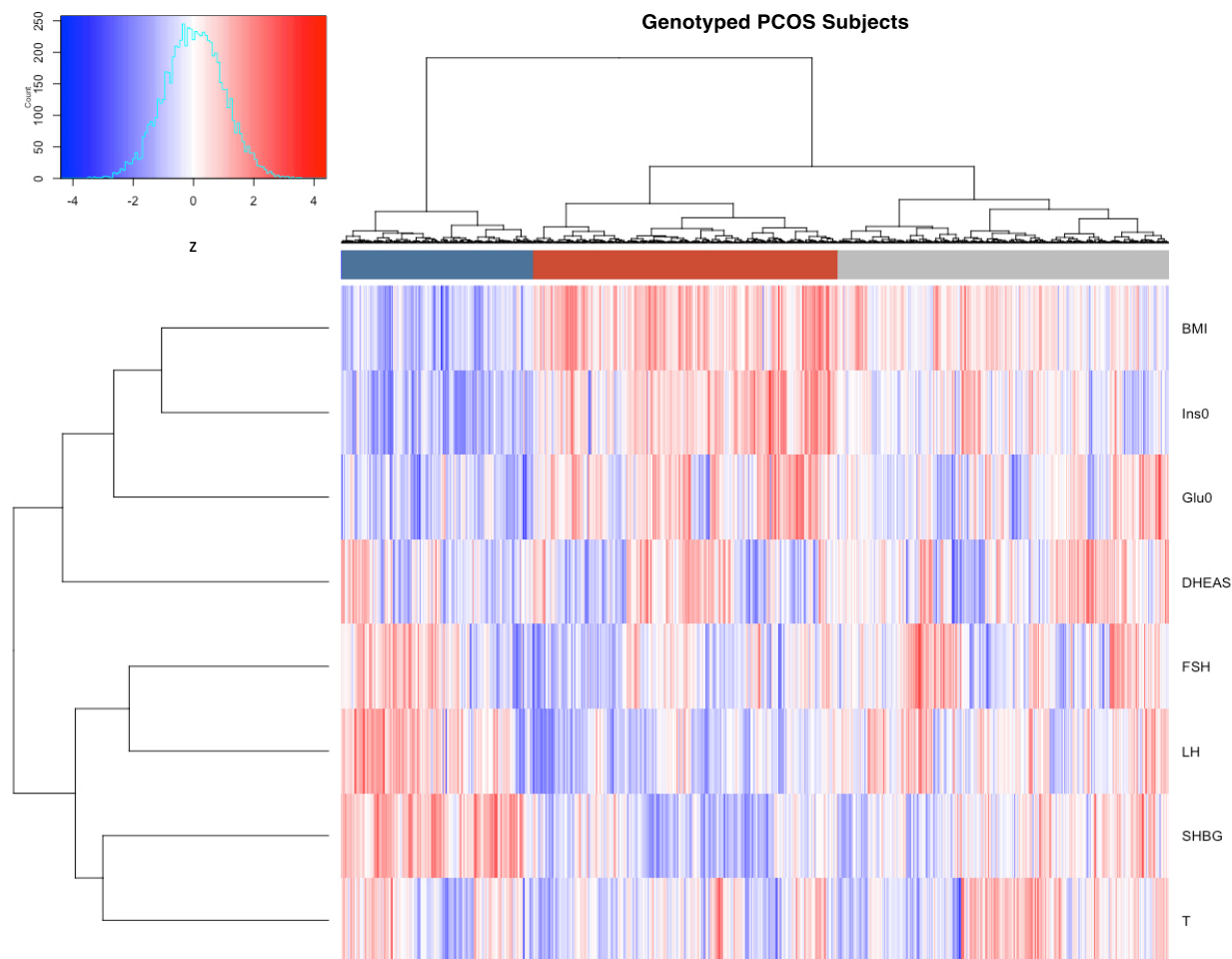
243 subtype with $\bar{\gamma}_C=0.55$. The indeterminate group did not appear to capture any discernable
244 pattern within the data ($\bar{\gamma}_C=0.41$) and was both overlapping and intermediate between the
245 reproductive and metabolic subtypes on the SHBG vs. Ins0/BMI axis.

246 The clustering procedure was then repeated in an independent, non-genotyped cohort
247 of 263 NIH PCOS cases diagnosed according to the same criteria as the genotyped cohort.

248 The clustering yielded similar results, with a comparable distribution of reproductive (26%,
249 $\bar{\gamma}_C=0.57$), metabolic (39%, $\bar{\gamma}_C=0.46$), and indeterminate clusters (35%, $\bar{\gamma}_C=0.40$) (**Fig 5**).

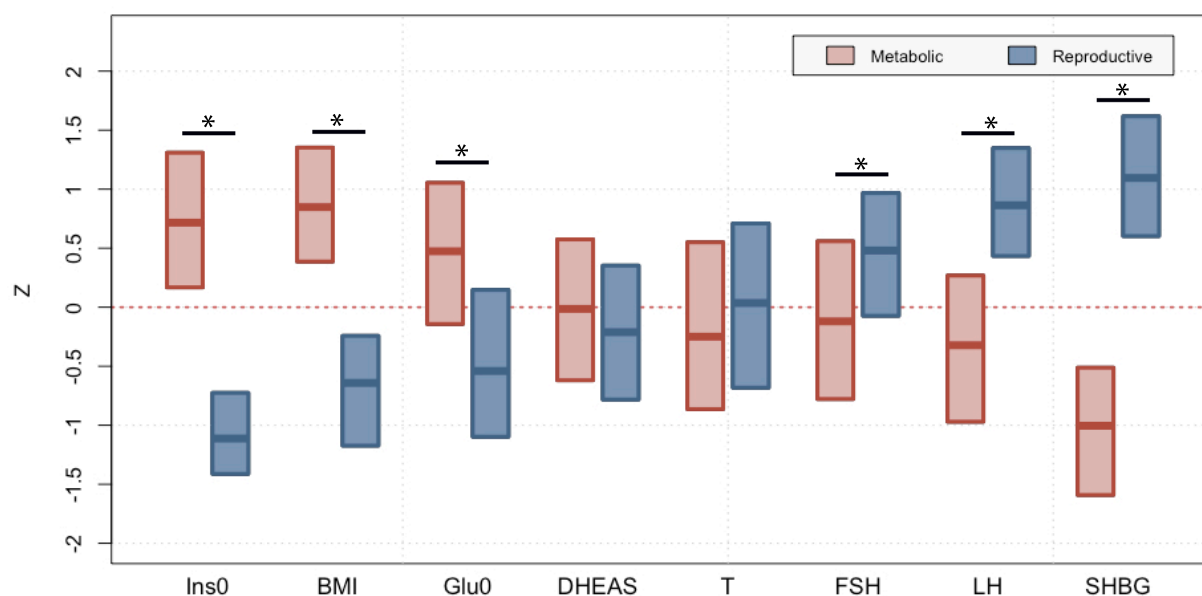
250

251



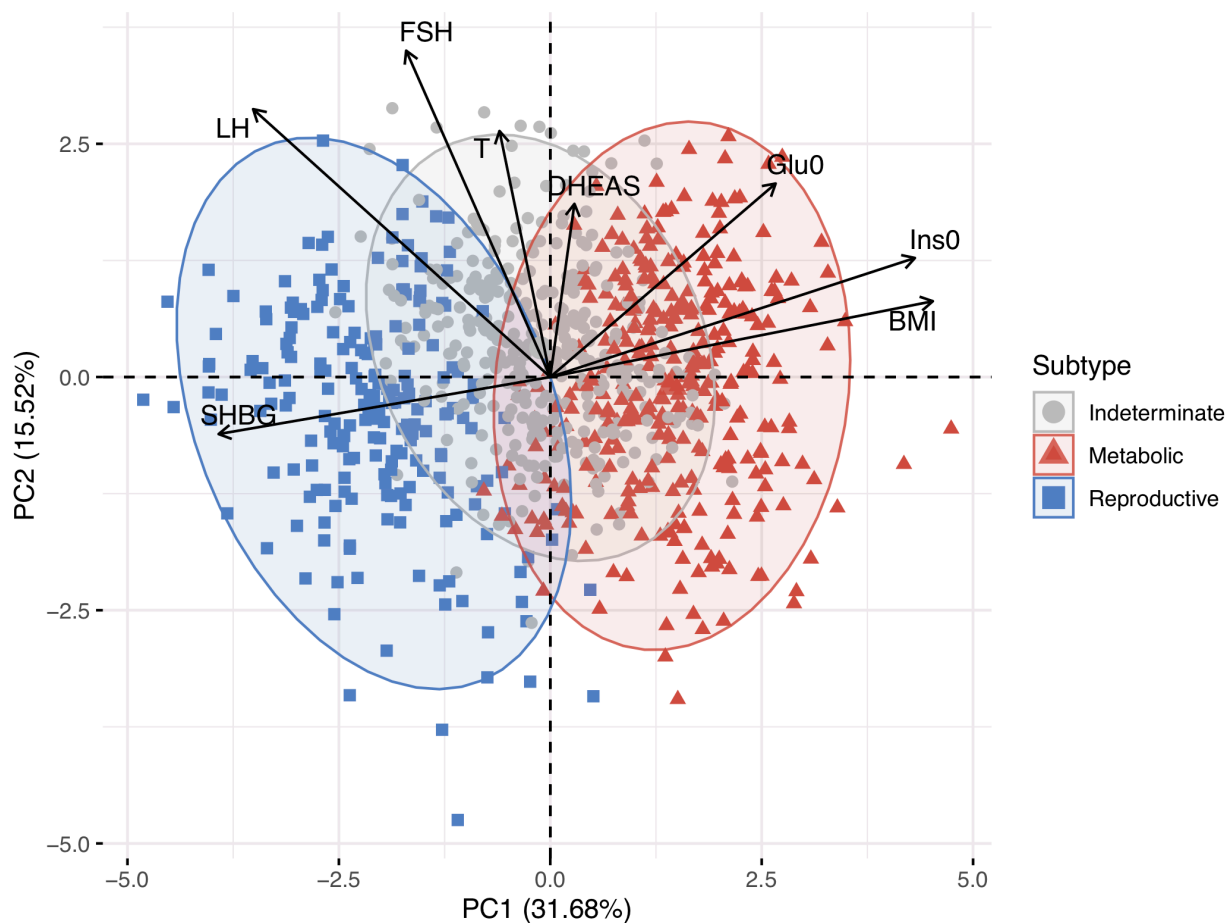
252

253 **Fig 2. Hierarchical clustering of genotyped PCOS cases.** Hierarchical clustering of 893
254 genotyped PCOS cases according to adjusted quantitative traits revealed two distinct
255 phenotypic subtypes: a “reproductive” cluster, and a “metabolic” cluster; the remaining
256 cases were designated as “indeterminate”. The reproductive, metabolic, and indeterminate
257 clusters are shown in the color bar as dark blue, dark red, and grey, respectively. Heatmap
258 colors correspond to trait z-scores, as shown in the frequency histogram where red
259 indicates high values and blue indicates low values for each trait. BMI, body mass index;
260 SHBG, sex hormone binding globulin; DHEAS, dehydroepiandrosterone sulfate; Glu0,
261 fasting glucose; Ins0, fasting insulin; LH, luteinizing hormone; FSH, follicle-stimulating
262 hormone.



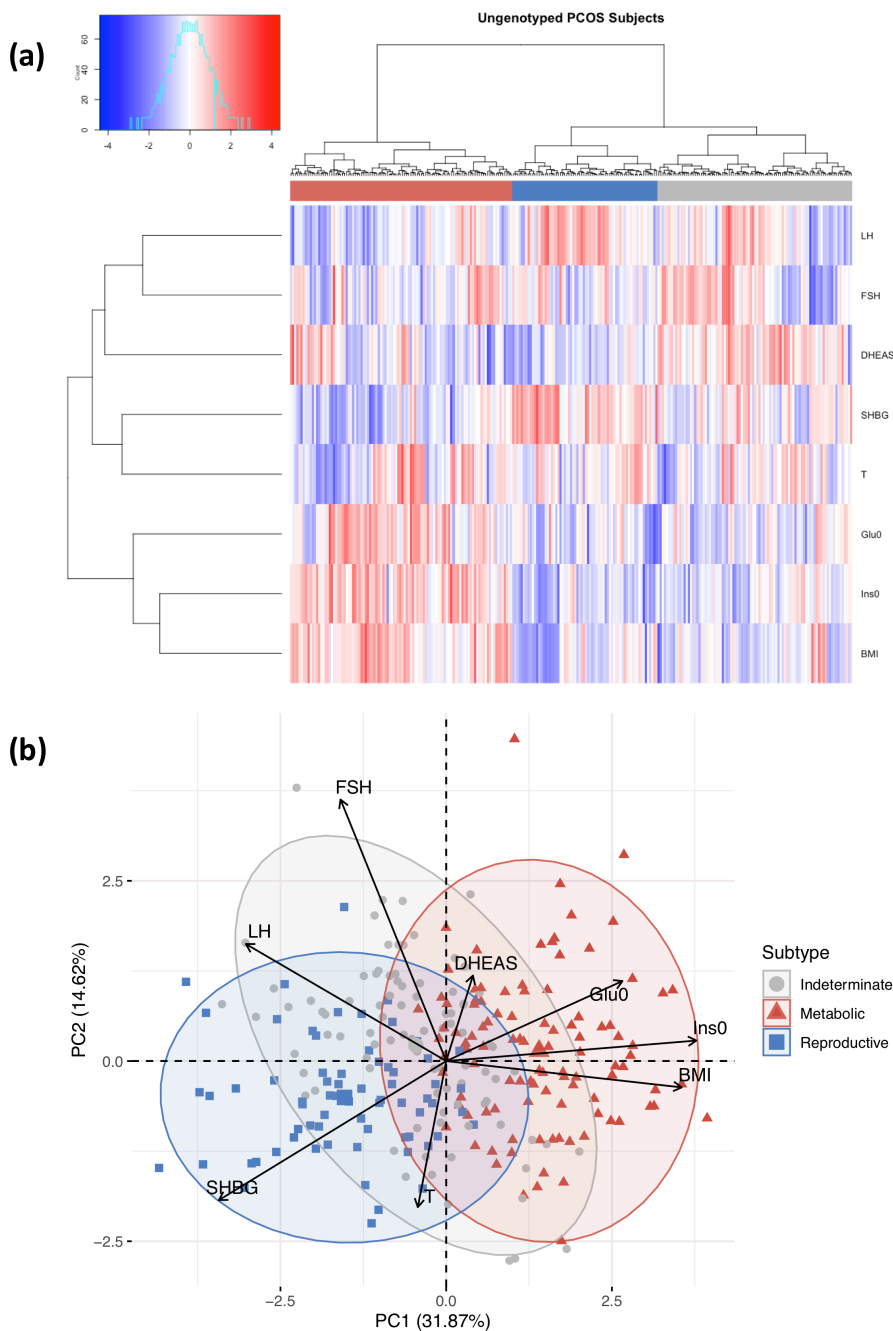
263

264 **Fig 3. Phenotypic trait distributions in reproductive and metabolic subtypes.** Median
265 and interquartile ranges are shown for normalized, adjusted quantitative trait distributions of
266 genotyped PCOS cases with reproductive or metabolic subtype. The figure illustrates the
267 traits for which the subtypes differ significantly with an asterisk (*Bonferroni adjusted
268 Wilcoxon, $P_{adj} < 0.05$): Ins0, BMI, Glu0, FSH, LH and SHBG. BMI, body mass index; SHBG,
269 sex hormone binding globulin; DHEAS, dehydroepiandrosterone sulfate; Glu0, fasting
270 glucose; Ins0, fasting insulin; LH, luteinizing hormone; FSH, follicle-stimulating hormone.



271

272 **Fig 4. PCA plot of quantitative traits for genotyped PCOS cases.** Genotyped PCOS
273 cases are plotted on the first two principal components of the quantitative trait data and
274 colored according to their identified subtype. Subtype clusters are shown as 95%
275 concentration ellipses, assuming bivariate normal distributions. The relative magnitude and
276 direction of trait correlations with the principal components are shown with black arrows.



277

278 **Fig 5. Clustering of ungenotyped PCOS cases.** (a) Hierarchical clustering of 263
279 ungenotyped PCOS cases according to adjusted quantitative traits replicate
280 reproductive (blue), metabolic (red), and unclassified (grey) clusters. Heatmap colors
281 correspond to trait z-scores. (b) PCA plot of ungenotyped PCOS cases replicate results
282 from genotyped cases. (a) Hierarchical clustering of 263 ungenotyped PCOS cases
283 according to adjusted quantitative traits replicate reproductive (blue), metabolic (red),
284 unclassified (grey) clusters. Heatmap colors correspond to trait z-scores. (b) PCA plot of
285 ungenotyped PCOS cases replicate results from genotyped cases.

286 Subtype genetic associations

287 Genome-wide association testing identified alleles in four novel loci that were
288 associated with the reproductive PCOS subtype at genome-wide significance (chr1 p36.21
289 *PRDM2/KAZN*, $P=2.23\times 10^{-10}$; chr2 q37.3 *IQCA1*, $P=2.76\times 10^{-9}$; chr4 q22.3
290 *BMPRI1/UNC5C*, $P=9.71\times 10^{-9}$; chr5 p14.2-p14.1 *CDH10*, $P=1.17\times 10^{-8}$) and one novel locus
291 that was significantly associated with the metabolic subtype (chr2 q24.2-q24.3
292 *KCNH7/FIGN*, $P=1.03\times 10^{-8}$). Association testing on the indeterminate cases replicated the
293 11p14.1 *FSHB* locus from our original GWAS (18) (**Table 1; Figs 6 and 7**).

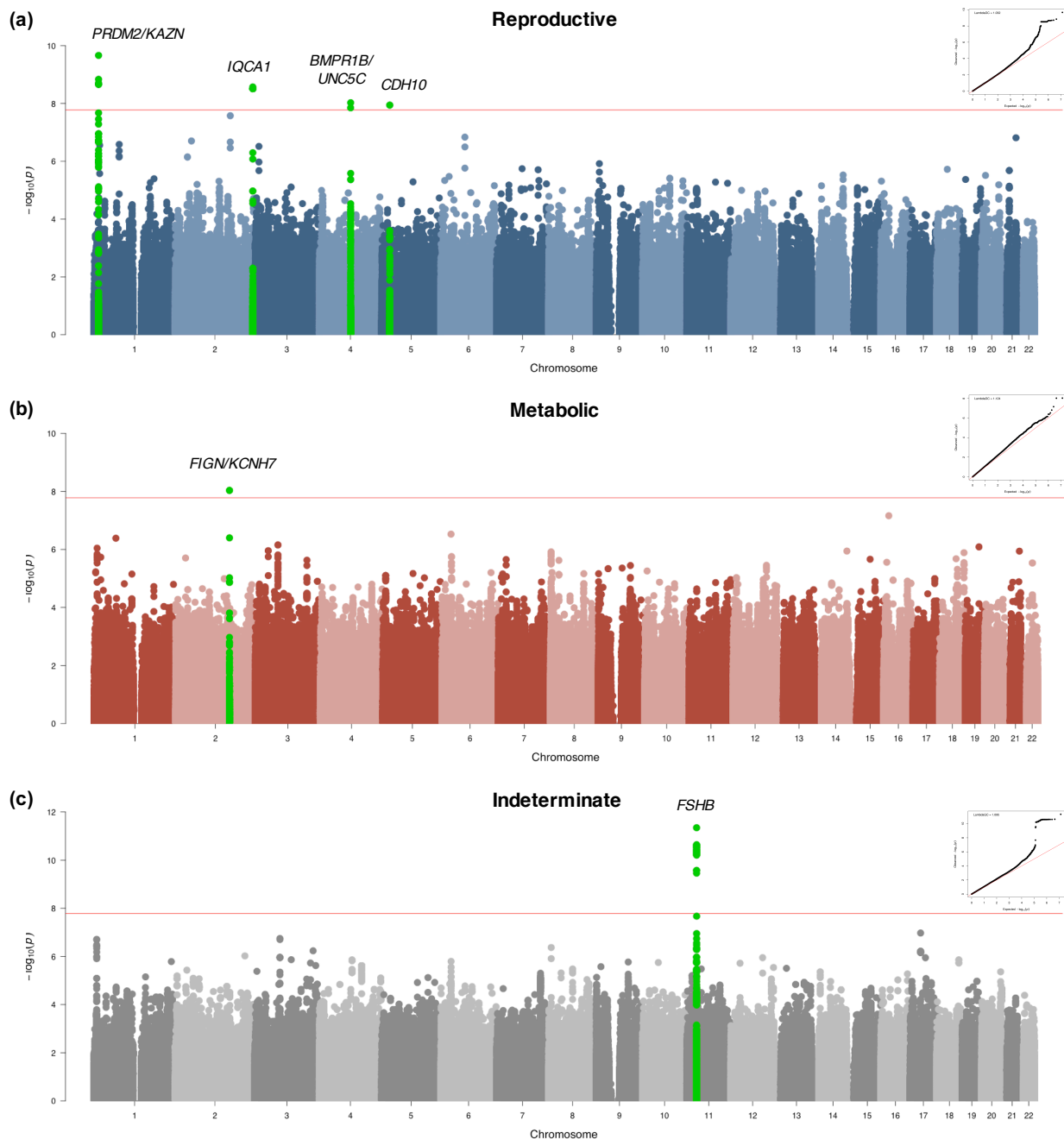
294 The strongest association signal with the reproductive subtype appeared in an
295 intergenic region of 1p36.21 579kb downstream of the *PRMD2* gene and 194kb upstream
296 from the *KAZN* gene (**Fig 8a**). The lead SNP in the locus (rs78025940; OR=4.75, 2.82-7.98
297 95%CI, $P_1=2.16\times 10^{-10}$, $P_{\text{meta}}=2.23\times 10^{-10}$) was imputed ($r^2 = 0.91$) in Stage 1 only. The SNP
298 was not genotyped in Stage 2. The lead genotyped SNP in the locus (rs16850259) was also
299 associated with the reproductive subtype with genome-wide significance ($P_{\text{meta}}=2.14\times 10^{-9}$)
300 and was genotyped only in Stage 1 (OR=5.57, 3.24-9.56 95% CI, $P_1=2.08\times 10^{-9}$). In ovarian
301 tissue, the SNPs appear to be centrally located within a large 2Mb TAD stretching from the
302 *FHADI* gene to upstream of the *PDPN* gene (**Fig 9**).

303 The 2q37.3 locus spanned a 50kb region of strong LD overlapping the 5' end and
304 promoter region of the *IQCA1* gene (**Fig 8b**). The SNP rs76182733 had the strongest
305 association in this locus ($P_{\text{meta}}=2.76\times 10^{-9}$) with the reproductive subtype. The signal was
306 genotyped only in Stage 1 (OR=5.68, 3.00-10.78 95%CI, $P_1=2.69\times 10^{-9}$) and was imputed
307 with an imputation r^2 value of 0.84.

308 **Table 1. Genome-wide significant associations with PCOS subtypes**

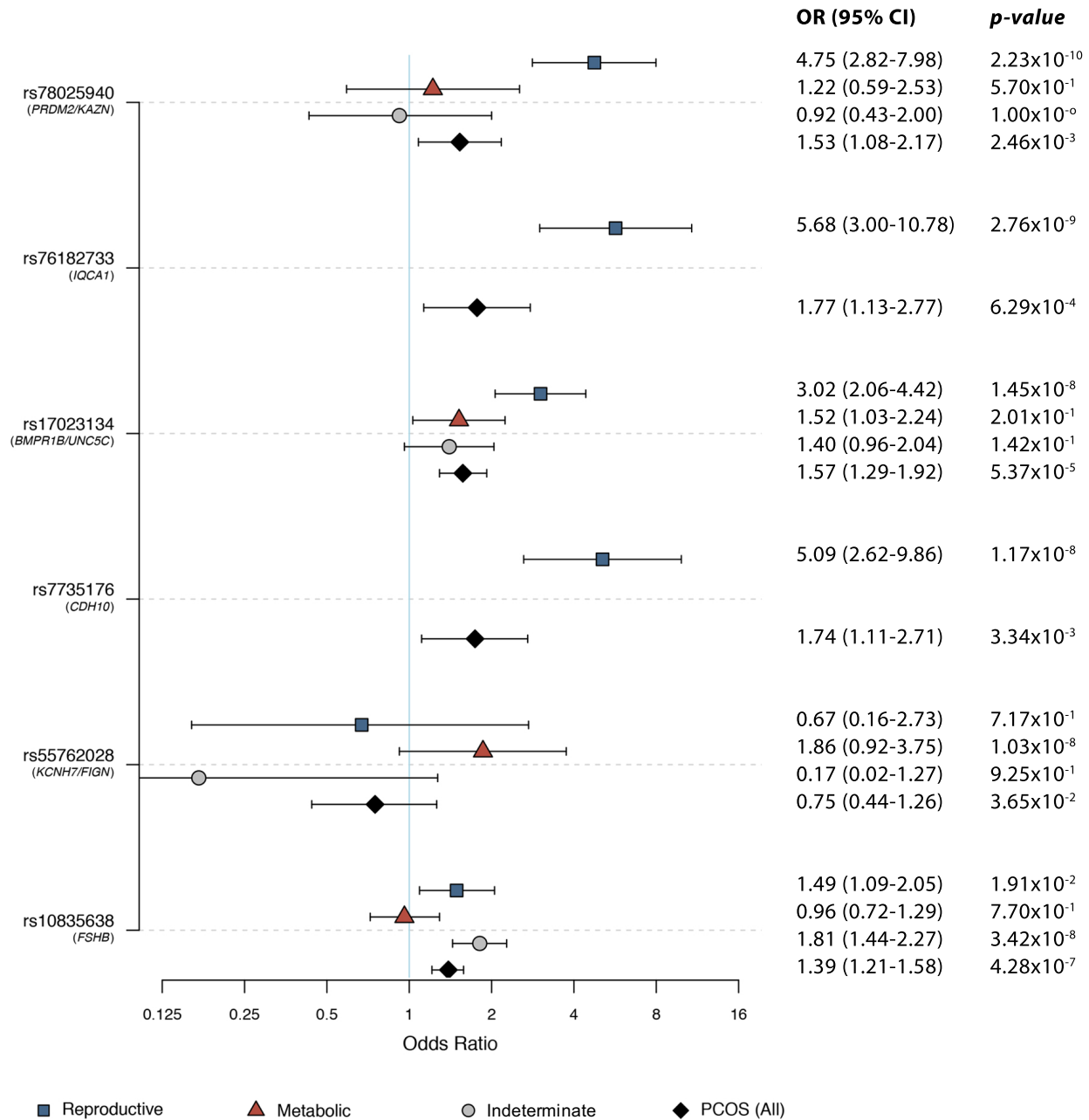
Chr	Mb	Variant	Gene(s)	EA	Stage 1 (Discovery)						Stage 2 (Replication)						P _{meta}
					EAF	β	OR	95% CI	P	Imp r ²	EAF	β	OR	95% CI	P	Imp r ²	
1	14.7	rs78025940	<i>PRDM2/ KAZN</i>	A	0.02	3.02	4.75	2.82-7.98	2.16 × 10 ⁻¹⁰	0.91	-	-	-	-	-	2.23 × 10 ⁻¹⁰	
2	237.4	rs76182733	<i>IQCA1</i>	G	0.01	3.79	5.68	3.00-10.78	2.67 × 10 ⁻⁹	0.84	-	-	-	-	-	2.76 × 10 ⁻⁹	
4	96.1	rs17023134	<i>BMPR1B/ UNC5C</i>	G	0.05	1.62	3.02	2.06-4.42	1.40 × 10 ⁻⁸	0.87	0.06	0.61	1.71	0.98-2.99	7.81 × 10 ⁻²	0.83	9.71 × 10 ⁻⁹
5	24.7	rs7735176	<i>CDH10</i>	A	0.01	3.80	5.09	2.62-9.86	1.14 × 10 ⁻⁸	0.93	-	-	-	-	-	1.17 × 10 ⁻⁸	
2	164.2	rs55762028	<i>KCNH7/ FIGN</i>	C	0.01	5.05	1.86	0.92-3.75	9.17 × 10 ⁻⁹	0.96	-	-	-	-	-	1.03 × 10 ⁻⁸	
11	30.3	rs10835638	<i>FSHB</i>	T	0.16	0.78	1.81	1.44-2.27	3.13 × 10 ⁻⁸	0.98	0.17	0.77	2.01	1.49-2.70	2.67 × 10 ⁻⁵	0.97	4.94 × 10 ⁻¹²

309 Variant information and association statistics are shown for the most strongly associated SNP in each significant locus.
 310 Reproductive subtype loci are highlighted in blue, metabolic loci in red, indeterminate loci in grey. EA = effect allele; EAF: effect
 311 allele frequency in cases and controls combined; β = effect size from association regression; OR = odds ratio; CI = confidence
 312 interval; Imp r² = imputation r² for imputed SNPs; P = stage-specific significance as assessed by logistic regression; P_{meta} =
 313 significance as assessed by sample-size weighted two-strata meta-analysis, adjusted for genomic inflation factor. Cases and
 314 controls by stage: Stage 1 = 201 metabolic, 123 lean, 231 indeterminate, 2964 controls; Stage 2 = 128 metabolic, 84 lean, 126
 315 indeterminate, 1134 controls. NOTE: Not all SNPs were genotyped or imputed in both stages.



316

317 **Fig 6. Genome-wide association results.** Manhattan plots for (a) reproductive and (b)
318 metabolic PCOS subtypes. The red horizontal line indicates genome-wide significance
319 ($p \leq 2.5 \times 10^{-8}$). Genome-wide significant loci are colored in green and labeled according
320 to nearby gene(s). QQ plots with genomic inflation factor, λ_{GC} , are shown adjacent to
321 corresponding plots.



322

323

324

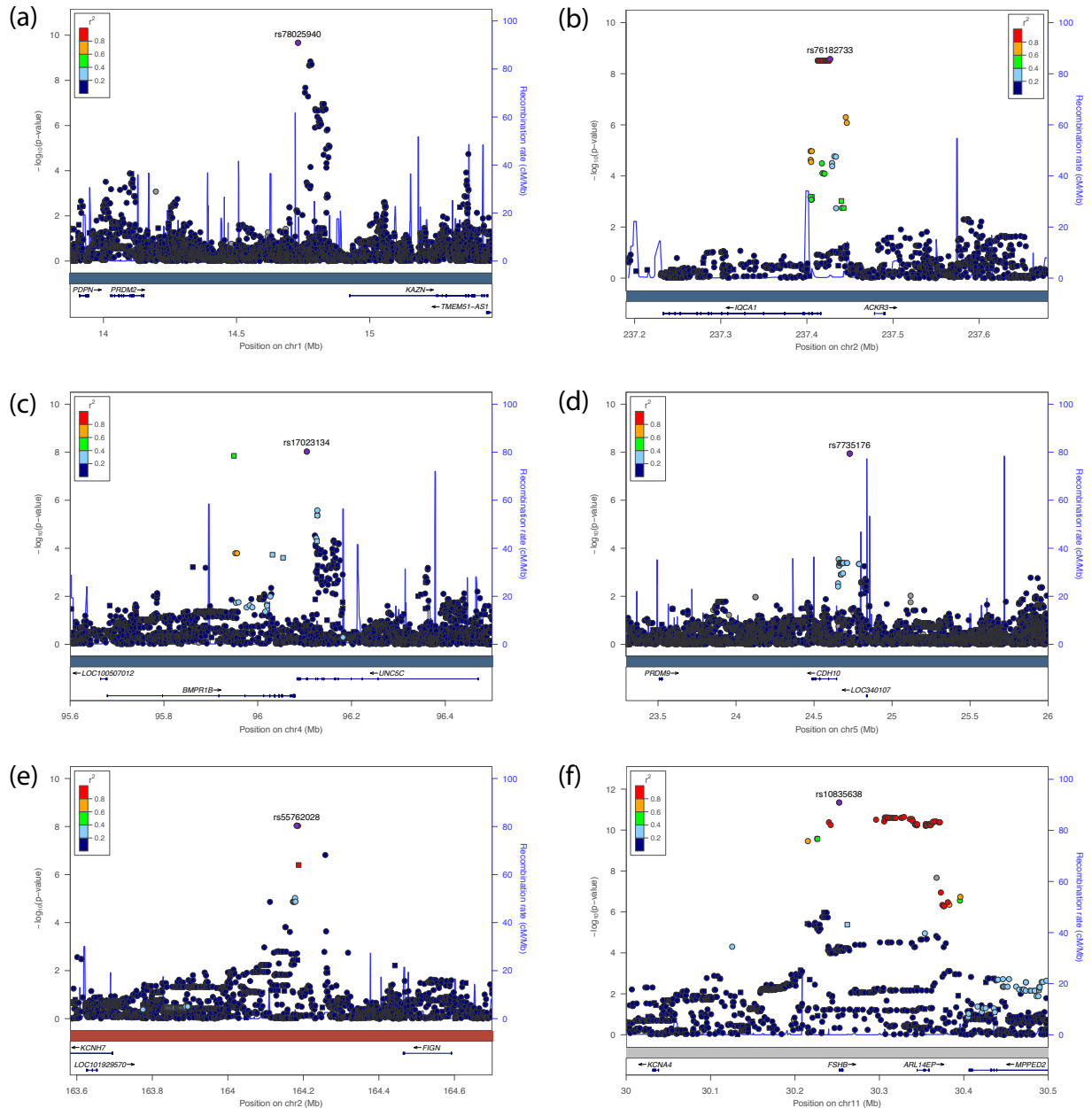
325

326

327

328

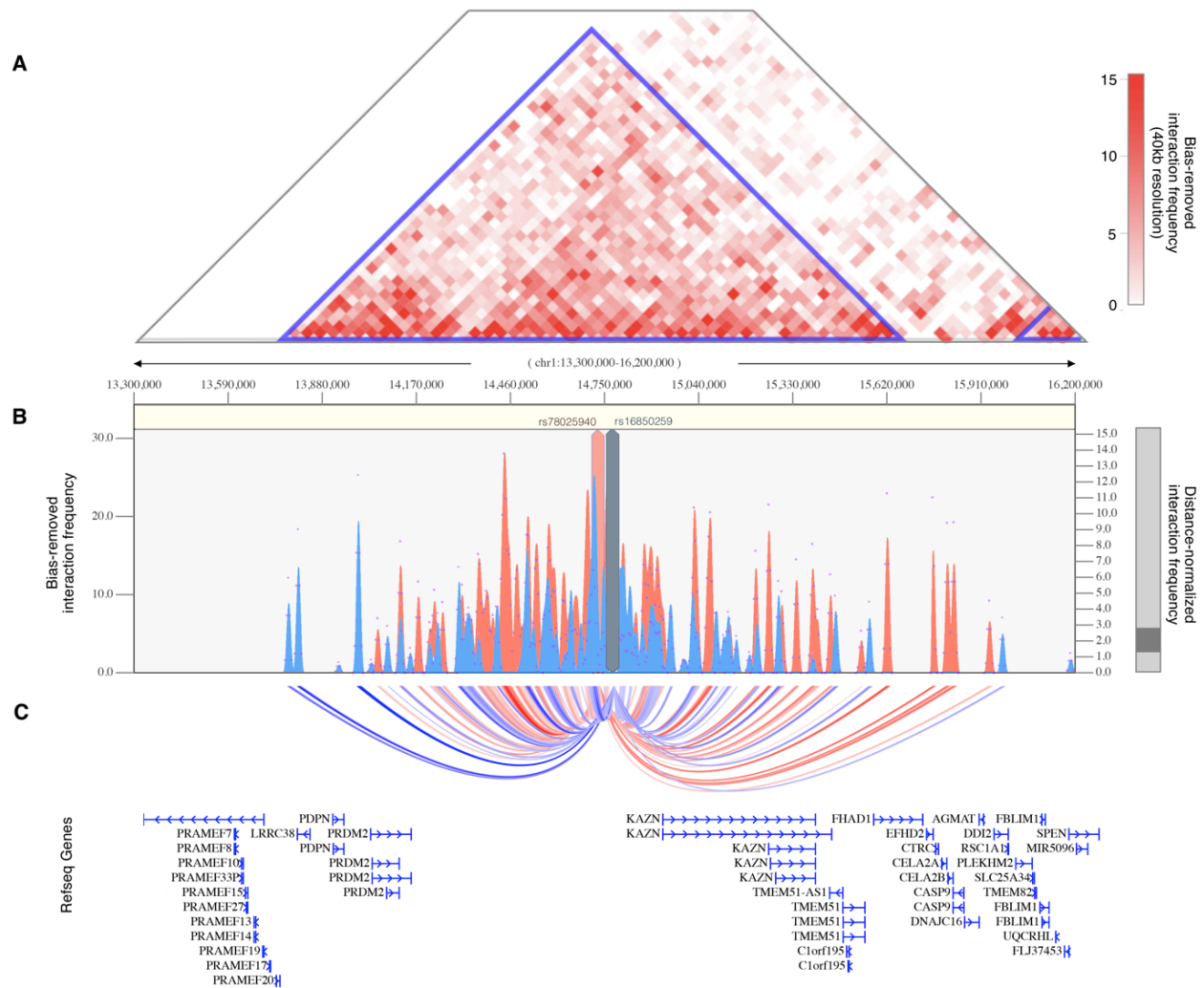
Fig 7. Risk allele odds ratios in PCOS and PCOS subtypes. Odds ratios with 95% confidence intervals and association P-values from the Stage 1 discovery cohort are shown for each subtype-specific novel risk allele identified in this study relative to the corresponding values for the other subtypes and for PCOS disease status in general (includes all subtypes). Some SNPs were not characterized in certain subtypes due to low allele counts or low imputation confidence.



329

330 **Fig 8. Regional association plots of genome-wide significant loci.** Regional plots of
 331 association (left y-axis) and recombination rates (right y-axis) for the chromosomes (a)
 332 1p36.21, (b) 2q37.3, (c) 4q22.3, (d) 5p14.2-p14.1, (e) 2p24.2-q24.3, and (f) 1p14.1 loci
 333 after imputation. The lead SNP in each locus is labeled and marked in purple. All other
 334 SNPs are color coded according to the strength of LD with the top SNP (as measured by
 335 r^2 in the European 1000 Genomes data). Imputed SNPs are plotted as circles and
 336 genotyped SNPs as squares.

337



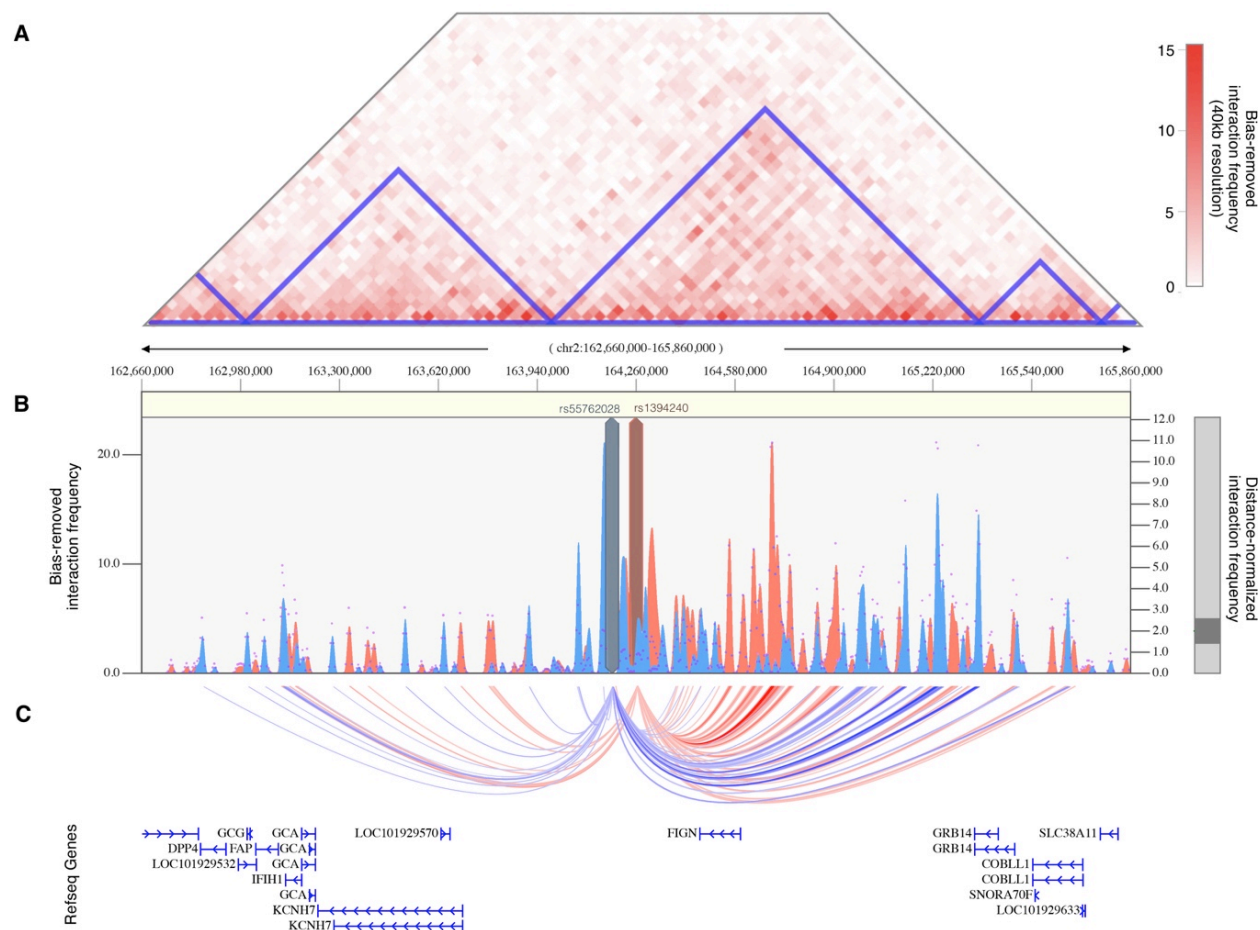
338

339 **Fig 9. Chromatin interaction map of *PRDM2/KAZN* locus.** (A) Shown is the
 340 interaction frequency heatmap from chr1:13,300,000-16,200,000 in ovarian tissue. The
 341 color of the heatmap indicates the level of normalized interaction frequencies with blue
 342 triangles indicating topological association domains. (B) One-to-all interaction plots are
 343 shown for the lead SNP (rs78025940; shown in red) and lead genotyped SNP
 344 (rs16850259; shown in blue) as bait. Y-axes on the left and the right measure bias-
 345 removed interaction frequency (red and blue bar graphs) and distance-normalized
 346 interaction frequency (magenta dots), respectively. (C) The arc-representation of
 347 significant interactions for distance-normalized interaction frequencies ≥ 2 is displayed
 348 relative to the RefSeq-annotated genes in the locus.

349 The 4q22.3 locus spanned a 500kb region of LD including the 3' ends of both the
350 *BMPRI1B* and *UNC5C* genes (**Fig 8c**). The most strongly associated SNP (rs17023134;
351 $P_{\text{meta}}=9.71\times 10^{-9}$) in the locus was within an intronic region of *UNC5C*, and was associated
352 with the reproductive subtype in the Stage 1 discovery cohort with genome-wide significance
353 (OR=3.02, 2.06-4.42 95%CI, $P_1=1.40\times 10^{-8}$), but was not significantly associated in the Stage
354 2 replication analysis (OR=1.71, 0.98-2.99 95%CI, $P_2=7.8\times 10^{-2}$). The SNP was imputed with
355 an r^2 of 0.87 and 0.83 in the Stage 1 and Stage 2 analyses, respectively. The most strongly
356 associated genotyped SNP in the locus (rs10516957) was also genome-wide significant
357 ($P_{\text{meta}}=1.46\times 10^{-8}$) and was located in an intronic region of *BMPRI1B*. The genotyped SNP was
358 nominally associated with the reproductive subtype in both the Stage 1 (OR=2.42, 1.66-3.52
359 95%CI, $P_1=6.72\times 10^{-6}$) and Stage 2 (OR=2.40, 1.51-3.82 95%CI, $P_2=4.7\times 10^{-4}$) analyses with
360 nearly identical effect sizes.

361 In the 5p14.2-p14.1 locus, 83kb upstream of the *CDH10* gene (**Fig 8d**), two adjacent
362 SNPs (rs7735176, rs16893866) in perfect LD were equally associated with the reproductive
363 subtype with genome-wide significance ($P_{\text{meta}}=1.17\times 10^{-8}$). The SNPs were imputed in Stage
364 1 (OR=5.09, 2.62-9.86 95%CI, $P_1=1.14\times 10^{-8}$) with an imputation r^2 of 0.93.

365 The single locus containing genome-wide significant associations with the metabolic
366 subtype was in an intergenic region of 2q24.2-q24.3 roughly 200kb downstream from *FIGN*
367 and 500kb upstream from *KCNH7* (**Fig 8e**). The lead SNP, rs55762028, was imputed in
368 Stage 1 only (OR=1.86, 0.92-3.75 95%CI, $P_1=9.17\times 10^{-9}$, $P_{\text{meta}}=1.03\times 10^{-8}$). In pancreatic
369 tissue, the lead SNPs appear to be located terminally within a 1.3Mb TAD encompassing the
370 *FIGN* gene and reaching upstream to the *GRB14* gene (**Fig 10**).



371

372 **Fig 10. Chromatin interaction map of *KCHN7/FIGN* locus.** (A) Shown is the
373 interaction frequency heatmap from chr2:162,660,000 to 165,860,000 in pancreatic
374 tissue. The color of the heatmap indicates the level of normalized interaction
375 frequencies with blue triangles indicating topological association domains. (B) One-to-all
376 interaction plots are shown for the lead SNP (rs13401392; shown in blue) and 2nd-
377 leading SNP (rs1394240; shown in red) as bait. Y-axes on the left and the right measure
378 bias-removed interaction frequency (blue and red bar graphs) and distance-normalized
379 interaction frequency (magenta dots), respectively. (C) The arc-representation of
380 significant interactions for distance-normalized interaction frequencies ≥ 2 is displayed
381 relative to the RefSeq-annotated genes in the locus.

382 Association testing on the indeterminate cases replicated the genome-wide significant
 383 association in the 11p14.1 *FSHB* locus (**Fig 8f**) identified in our original GWAS (14). The
 384 lead SNP (rs10835638; $P_{\text{meta}}=4.94 \times 10^{-12}$) and lead genotyped SNP (rs10835646;
 385 $P_{\text{meta}}=2.75 \times 10^{-11}$) in this locus differed from the index SNPs identified in our original GWAS
 386 (rs11031006) and in the PCOS meta-analysis (rs11031005), but both of the previously
 387 identified index SNPs were also associated with the indeterminate subgroup with genome-
 388 wide significance in this study (rs11031006: $P_{\text{meta}}=2.96 \times 10^{-10}$; rs11031005: $P_{\text{meta}}=2.91 \times 10^{-10}$)
 389 and are in LD with the lead SNP rs10835638 ($r^2 = 0.59$) (38). The other significant signals
 390 from our original GWAS (18) were not reproduced in any of the subtype association tests
 391 performed in this study (**Table 2**).

392

393 **Table 2. Previous GWAS association signals in PCOS subtypes**

Variant	Locus	PCOS	Reproductive	Metabolic	Indeterminate
rs804279	<i>GATA4/NEIL2</i>	$P = 8.0 \times 10^{-10}$	$P = 2.4 \times 10^{-3}$	$P = 9.9 \times 10^{-2}$	$P = 3.1 \times 10^{-3}$
rs10993397	<i>C9orf3</i>	$P = 4.6 \times 10^{-13}$	$P = 2.3 \times 10^{-4}$	$P = 6.9 \times 10^{-5}$	$P = 1.1 \times 10^{-5}$
rs11031006	<i>FSHB</i>	$P = 1.9 \times 10^{-8}$	$P = 8.8 \times 10^{-6}$	$P = 6.6 \times 10^{-1}$	$P = 3.0 \times 10^{-10}$

394 Subtype-specific association statistics are shown for each of the SNPs that were
 395 significantly associated with PCOS in Hayes and Urbanek et al. (18). P = significance as
 396 assessed by sample-size weighted two-strata meta-analysis, adjusted for genomic inflation.

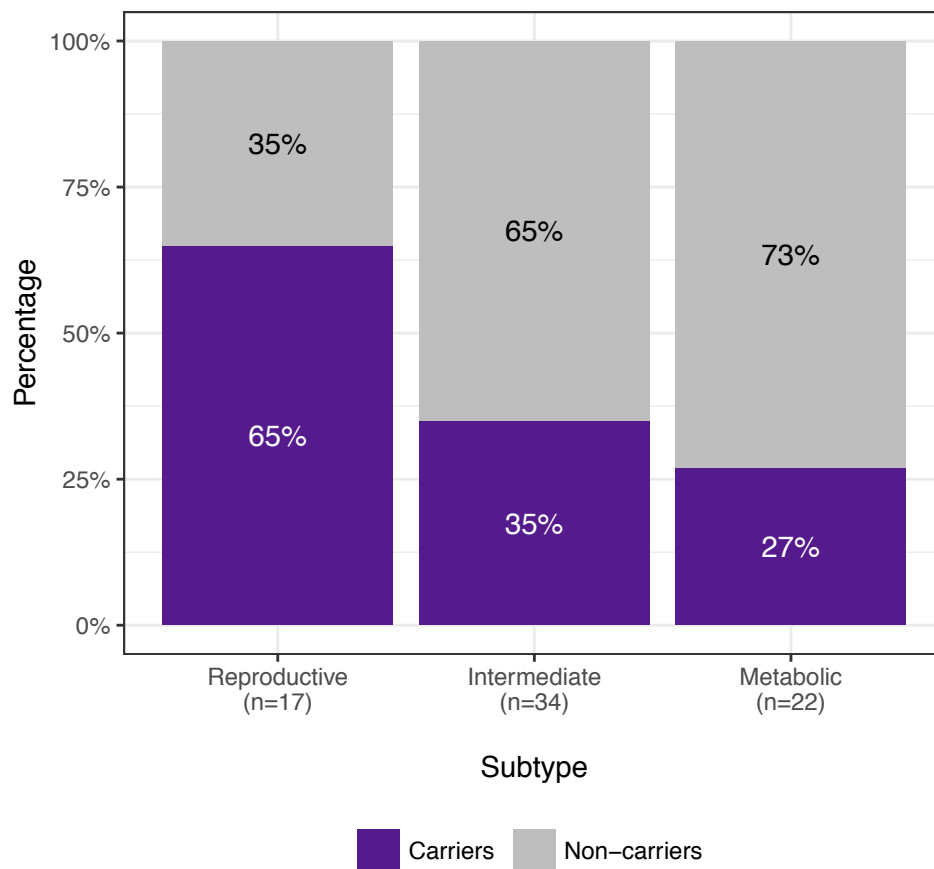
397

398 Subtypes in PCOS families

399 The RF classifier yielded the lowest overall subtype misclassification rate (13.2%) of
 400 the tested methods, according to 10-fold cross-validation of the genotyped cohort. Affected
 401 women from the family-based cohort were classified accordingly using a RF model. Seventy-
 402 three daughters of the 83 affected women from the family-based cohort had complete

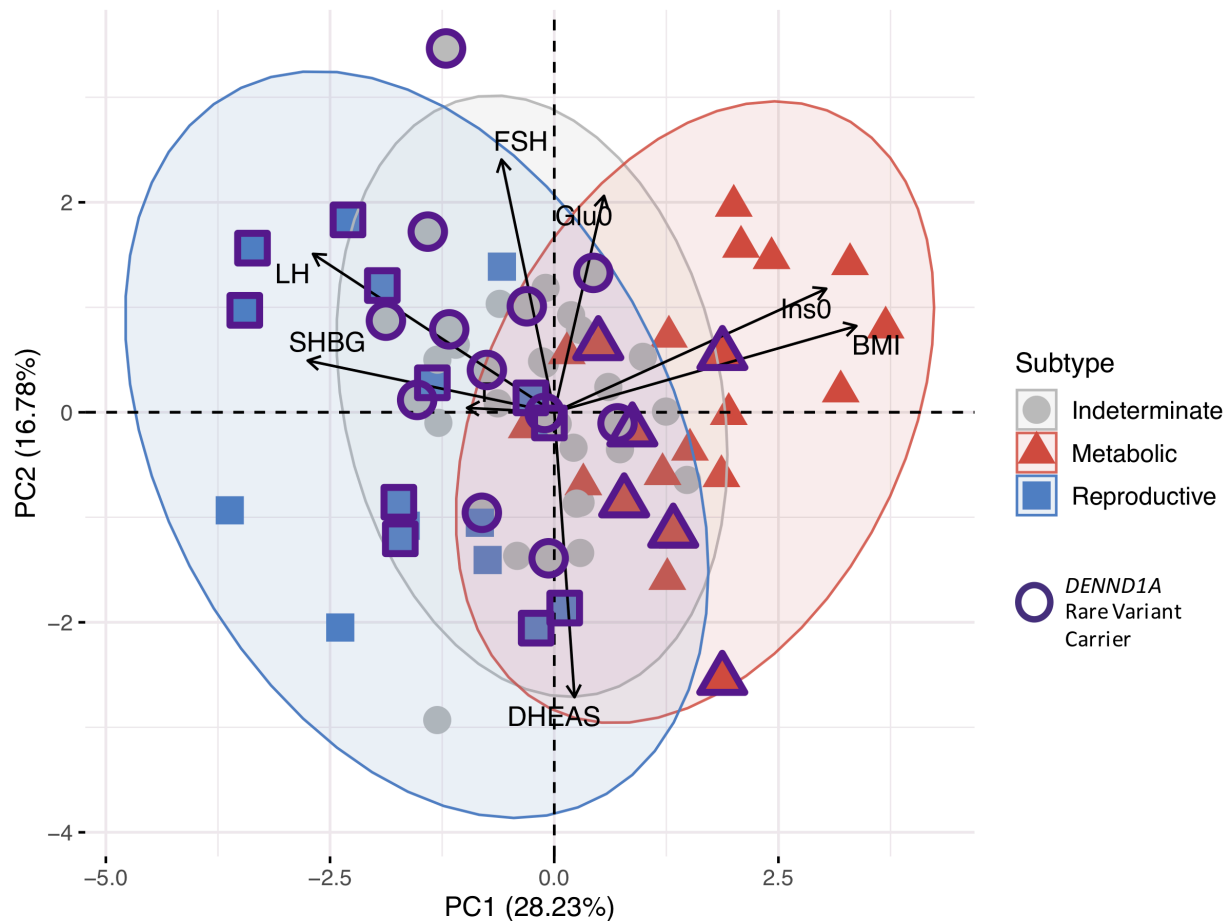
403 quantitative trait data available for subtype classification. Seventeen (23.3%) were classified
404 as having the reproductive subtype of PCOS, and 22 (30.1%) were classified as having the
405 metabolic subtype. Of 14 subtyped sibling pairs, only 8 were concordantly classified
406 (57.1%); however, there was only one instance of the reproductive subtype and metabolic
407 subtype occurring within the same nuclear family as the remaining discordant pairs each
408 featured one indeterminate member. The proportion of affected women with one or more of
409 the previously-identified (29) deleterious, rare variants in *DENNDIA* varied by subtype.
410 Women classified as having the reproductive subtype of PCOS were significantly more likely
411 to carry one or more of the *DENNDIA* rare variants compared to other women with PCOS
412 ($P=0.03$; **Fig 11**). The distribution of affected women and *DENNDIA* rare variant carriers are
413 shown relative to the adjusted quantitative trait PCs in **Fig 12**.

414



415

416 **Fig 11. *DENND1A* rare variant carriers by subtype.** The proportions of affected
417 women with *DENND1A* rare variants in families with PCOS are shown by classified
418 subtype. Women with the reproductive subtype were significantly more likely to
419 carry one or more of the *DENND1A* rare variants compared to other women with
420 PCOS (P=0.03)



421

422 **Fig 12. PCA of affected women in PCOS families showing *DENND1A* rare variant**
423 **carriers.** The distribution of affected women and *DENND1A* rare variant carriers are
424 shown relative to their classified subtypes and their adjusted quantitative trait PCs in
425 families affected by PCOS.

426

427 Discussion

428 It is becoming increasingly clear that common, complex traits, such as T2D, are a
429 heterogeneous collection of disease subtypes (24,45-47). There is emerging evidence that these
430 subtypes have different genetic architectures (24,47,48). Consistent with these concepts, we
431 identified reproductive and metabolic subtypes of PCOS by unsupervised hierarchical cluster
432 analysis of quantitative hormonal traits and BMI and found novel loci uniquely associated with

433 these subtypes with substantially larger effect sizes than those associated with PCOS disease
434 status in GWAS (18-22). These findings suggest that these subtypes are both genetically distinct
435 as well as more etiologically homogenous (9). Our findings are in contrast to the recent PCOS
436 GWAS meta-analysis (21) that found that only one of 14 loci was uniquely associated with the
437 NIH phenotype compared to non-NIH Rotterdam phenotypes. These latter findings suggest that
438 the NIH and Rotterdam diagnostic criteria do not identify biologically distinct subtypes of
439 PCOS. There have been previous efforts to subtype PCOS using unsupervised clustering (25-28),
440 but no subsequent investigation into the biologic relevance of the resulting subtypes.

441 The key traits driving the subtypes were BMI, insulin, SHBG, glucose, LH, and FSH
442 levels. The reproductive subtype was characterized by higher LH and SHBG levels with
443 lower BMI and blood glucose and insulin levels. The metabolic subtype was characterized by
444 higher BMI and glucose and insulin levels with relatively low SHBG and LH levels. The
445 remaining 40% of cases had no distinguishable cluster-wide characteristics and mean trait
446 values were between those of the reproductive and metabolic subtypes. The relative trait
447 distributions and results of the PCAs (**Figs 3, 4, 5b**) showed the reproductive and metabolic
448 subtypes as collections of subjects on opposite ends of a phenotypic spectrum with the
449 remaining indeterminate subjects scattered between the two. Bootstrapping and clustering in
450 an independent cohort revealed that the reproductive and metabolic subtypes were stable and
451 reproducible. When the GWAS was repeated, novel susceptibility loci were associated with
452 the reproductive and metabolic subtypes, suggesting that they had distinct genetic
453 architecture. The indeterminate PCOS cases were associated with the reported locus at
454 *FSHB*, but the association signal was stronger than that of our original GWAS (18),

455 suggesting that the indeterminate group was also more genetically homogenous after the
456 reproductive and metabolic subtypes were removed from the analysis.

457 Two of the loci associated with the reproductive subtype implicate novel biologic
458 pathways in PCOS pathogenesis. The association signal on chr1 appeared downstream of and
459 within the same TAD as the *PRDM2* gene (**Figs 8a, 9**), which is an estrogen receptor co-
460 activator (49) that is highly expressed in the ovary (50) and pituitary gland (51). In an
461 independent rare variant association study in PCOS families, *PRDM2* demonstrated the 5th
462 strongest gene-level association with altered hormonal levels in PCOS families ($P=6.92\times 10^{-3}$)
463 out of 339 genes tested (29). *PRDM2* binds with ligand bound estrogen receptor alpha ($ER\alpha$)
464 to open chromatin at $ER\alpha$ target genes (49,52). *PRDM2* also binds with the retinoblastoma
465 protein (53), which has been found to play an important role in follicular development in
466 granulosa cells (54,55).

467 The reproductive subtype association in the 4q22.3 locus overlapped with the
468 *BMPR1B* gene, which transcribes a type-I AMH receptor highly expressed in granulosa cells
469 and in GnRH neurons (16) that regulates follicular development (56). *BMPR1B* (bone
470 morphogenetic protein receptor type IB), also known as ALK6 (Activin Receptor-Like
471 Kinase 6), heterodimerizes with the TGF- β type-II receptors, including AMHR2, and binds
472 AMH and other BMP ligands to initialize TGF- β signaling via Smads 1/5/8 (57). *BMPR1B*
473 has been found to mediate the AMH response in ovine granulosa cells (58), and *BMPR1B*-
474 deficient mice are infertile and suffer from a variety of functional defects in the
475 ovary (59,60). One of the *BMPR1B* ligand genes, *BMP6*, had the 3rd strongest gene-level
476 association with altered hormonal levels ($P=4.00\times 10^{-3}$) out of 339 genes tested in our rare
477 variant association study in PCOS families (29). Collectively, these results make *BMPR1B* a

478 compelling candidate gene in PCOS pathogenesis. These findings also support our
479 sequencing studies that have implicated pathogenic variants in the AMH signaling pathway
480 in PCOS (61,62).

481 The nature of the potential involvement in PCOS is less clear for the other loci
482 associated with the reproductive subtype. The 2q37.3 locus overlapped with the promoter
483 region of the *IQCA1* gene. Its function in humans is not well characterized, but *IQCA1* is
484 highly expressed in the pituitary gland (51). The 5p14.2-p14.1 locus overlapped the promoter
485 region of the *CDH10* gene (Cadherin 10). *CDH10* is almost exclusively expressed in the
486 brain (50), and is putatively involved in synaptic adhesions, axon outgrowth and
487 guidance (63).

488 The lone significant association signal with the metabolic subtype was located in an
489 intergenic region 200-280kb downstream of the *FIGN* gene, 490-570kb upstream of *KCNH7*.
490 *KCNH7* encodes a voltage-gated potassium channel (alias ERG3). *KCNH7* is primarily
491 expressed in the nervous system (64), but has been found in murine islet cells (65,66). *FIGN*
492 encodes fidgetin, a microtubule-severing enzyme most highly expressed in the pituitary
493 gland and ovary (50). A genetic variant in *FIGN* was found to reduce the risk of congenital
494 heart disease in Han Chinese by modulating transmembrane folate transport (67,68). The
495 TAD encompassing the association signal in this locus includes *FIGN* and extends upstream
496 to the *GRB14* gene (**Fig 9**). *GRB14* plays an important role in insulin receptor
497 signaling (69,70) and has been associated with T2D in GWAS (71). Given the various
498 metabolic associations for the genes in this chromosomal region, it is plausible that causal
499 variants in this locus could impact a combination of these genes.

500 Despite evidence linking neighboring genes to PCOS pathways in each of the
501 aforementioned loci, it remains possible, of course, that other more distant genes in LD
502 underlie the association signals. Causal variants are often up to 2 Mb away from the
503 associated SNP, not necessarily in the closest gene (72). Fine-mapping and functional studies
504 are needed in order to confirm the causal variants in each of these loci. In addition, the
505 sample sizes for the subtype GWAS were small, some of the associations were based only on
506 imputed SNPs in Stage 1, and a replication association study has not yet been performed.
507 However, the aforementioned functional evidence for several of the loci—particularly for
508 *PRDM2* and *BMPRI1B*—support the validity of their associations. Also, the fact that one of
509 the genes associated with the reproductive subtype, *PRDM2*, was associated with PCOS
510 quantitative traits in our family-based analysis (29) does represent a replication of this signal
511 by an independent analytical approach. Nevertheless, our genetic association results should
512 be considered preliminary.

513 The effect sizes of the subtype alleles, particularly those associated with the
514 reproductive subtype (Odds Ratio [OR] 3.02-5.68) (**Table 1**), were substantially greater than
515 the effects (OR 0.70-1.51) observed for alleles associated with PCOS diagnosis in previous
516 GWAS (18-22). In general, there is an inverse relationship between allele frequency and
517 effect size (1) because alleles with larger phenotypic effects are subject to purifying selection
518 and, therefore, occur less frequently in the population (73,74). Accordingly, in contrast to the
519 common variants (Effect Allele Frequency [EAF]>0.05) associated with PCOS in previous
520 GWAS (18-22), the alleles associated with the subtypes were all of low frequency (EAF
521 0.01-0.05; **Table 1**). However, given the limited cohort sizes in this study, the subtype
522 association testing did not have adequate power to detect associations with more modest

523 effect sizes, such as those from our previous GWAS (18). It is also possible that the large
524 effect sizes were somewhat inflated by the so-called “winners curse” (75,76), but they
525 nonetheless suggest that the subtypes were more genetically homogeneous than PCOS
526 diagnosis in general.

527 In applying a subtype classifier to our family-based cohort, we found twelve affected
528 sibling pairs in which at least one of the daughters was classified with the reproductive or
529 metabolic subtype. Six of these sibling pairs were classified with the same subtype. There
530 was only one discordant pairing of the reproductive subtype with the metabolic subtype. This
531 further suggests that the reproductive and metabolic subtypes are genetically distinct in their
532 origins. The greater prevalence of *DENNDIA* rare variant carriers observed in women with
533 the reproductive subtype in the family-based cohort implicates this gene in the pathogenesis
534 of this subtype. *DENNDIA* is known to regulate androgen biosynthesis in the ovary (77,78);
535 therefore, we would expect *DENNDIA*-mediated PCOS to be more closely associated with
536 the reproductive subtype of PCOS. However, we did not find an association between any
537 *DENNDIA* alleles and the reproductive subtype in the subtype GWAS, perhaps due to allelic
538 heterogeneity or to our limited power to detect associations with more modest effect sizes.

539 We only studied women with PCOS as defined by the NIH diagnostic criteria. Future
540 studies will investigate whether similar reproductive and metabolic clusters are present in
541 non-NIH Rotterdam PCOS cases. In particular, it is possible that there will be no metabolic
542 subtype in non-NIH Rotterdam PCOS cases since these phenotypes have minimal metabolic
543 risk (79,80). Indeed, in a previous effort to identify phenotypic subtypes in Rotterdam PCOS
544 cases (28), the cluster that most closely resembled the reproductive subtype represented the
545 largest proportion of PCOS women at 44%, of whom only 78% met the NIH criteria for PCOS,

546 whereas the cluster that most closely resembled the metabolic subtype constituted only 12% of
547 the cohort, but 98% met the NIH diagnostic criteria. Due to the within-cohort normalization of
548 quantitative traits prior to clustering, our method is well-suited for identifying subsets of cases
549 that occupy either end of a phenotypic spectrum within different populations, but it therefore
550 may not be suitable for directly comparing subtype membership between different populations.

551 Our cohort included only women of European ancestry. It will be of considerable
552 importance to investigate whether subtypes are present in women with PCOS of other ancestries.
553 Women with PCOS of diverse races and ethnicities have similar reproductive and metabolic
554 features (81-83). However, there are differences in the severity of the metabolic defects due to
555 differences in the prevalence of obesity (84) and well as to racial/ethnic differences in insulin
556 sensitivity (85,86). Further, the susceptibility loci associated with subtypes in other ancestry
557 groups may differ since the low frequency and large effect size of the variants associated with
558 the reproductive subtype in our European cohort suggests these variants are of relatively recent
559 origin, and, therefore, may be population-specific (1,87,88).

560 In conclusion, using an unsupervised clustering approach featuring quantitative
561 hormonal and anthropometric data, we identified novel reproductive and metabolic subtypes
562 of PCOS with distinct genetic architectures. The genomic loci that were significantly
563 associated with either of these subtypes include a number of new, highly plausible PCOS
564 candidate genes. Moreover, our results demonstrate that precise phenotypic delineation,
565 resulting in more homogeneous subsets of affected individuals, can be more powerful and
566 informative than increases in sample size for genetic association studies. Our findings
567 indicate that further study into the genetic heterogeneity of PCOS is warranted and could
568 lead to a transformation in the way PCOS is classified, studied, and treated.

569 **Acknowledgements**

570 This study was supported by National Institutes of Health (NIH) Grants
571 R01HD057223 (A.D.), P50 HD044405 (A.D.) and R01 HD085227 (A.D.). M.D. was
572 supported by a Ruth L. Kirschstein National Research Service Award Institutional Research
573 Training Grant, T32 DK007169. We thank the NIH Cooperative Multicenter Reproductive
574 Medicine Network (<https://www.nichd.nih.gov/research/supported/rmn>) for recruiting some
575 of the women with PCOS who participated who participated in the genomewide association
576 study of Hayes and Urbanek et al. (18) and whose genotype data were used in this study . We
577 also thank the following investigators for recruiting some of the control women who
578 participated in the genomewide association study of Hayes and Urbanek et al. (18) and
579 whose genotype data were used in this study: Dimitrios Panidis, MD, PHD (Aristotle
580 University of Thessaloniki, Greece); Mark O. Goodarzi (Cedars-Sinai Medical Center, Los
581 Angeles, CA); Corrine K. Welt, MD (University of Utah School of Medicine, Salt Lake City,
582 UT; formerly of Massachusetts General Hospital, Boston, MA); Ahmed H. Kissebah
583 (deceased, Medical College of Wisconsin, Milwaukee, WI); Ricardo Azziz, MD (State
584 University of New York, NY; formerly of University of Alabama at Birmingham, AL); and
585 Evanthia Diamanti-Kandarakis, MD, PhD (University of Athens Medical School, Greece).

References

1. Manolio TA, Collins FS, Cox NJ, Goldstein DB, Hindorff LA, Hunter DJ, et al. Finding the missing heritability of complex diseases. *Nature*. 2009;461(7265):747-53.
2. Boyle EA, Li YI, Pritchard JK. An Expanded View of Complex Traits: From Polygenic to Omnigenic. *Cell*. 2017;169(7):1177-86.
3. Wray NR, Wijmenga C, Sullivan PF, Yang J, Visscher PM. Common Disease Is More Complex Than Implied by the Core Gene Omnigenic Model. *Cell*. 2018;173(7):1573-80.
4. Ringman JM, Goate A, Masters CL, Cairns NJ, Danek A, Graff-Radford N, et al. Genetic heterogeneity in Alzheimer disease and implications for treatment strategies. *Curr Neurol Neurosci Rep*. 2014;14(11):499.
5. Flint J, Kendler KS. The genetics of major depression. *Neuron*. 2014;81(3):484-503.
6. von Coelln R, Shulman LM. Clinical subtypes and genetic heterogeneity: of lumping and splitting in Parkinson disease. *Curr Opin Neurol*. 2016;29(6):727-34.
7. Ahlqvist E, Storm P, Karajamaki A, Martinell M, Dorkhan M, Carlsson A, et al. Novel subgroups of adult-onset diabetes and their association with outcomes: a data-driven cluster analysis of six variables. *Lancet Diabetes Endocrinol*. 2018;6(5):361-9.
8. Udler MS, Kim J, von Grotthuss M, Bonas-Guarch S, Cole JB, Chiou J, et al. Type 2 diabetes genetic loci informed by multi-trait associations point to disease mechanisms and subtypes: A soft clustering analysis. *PLoS Med*. 2018;15(9):e1002654.
9. Saria S, Goldenberg A. Subtyping: What it is and its role in precision medicine. *IEEE Intelligent Systems*. 2015;30(4):70-5.
10. Diamanti-Kandarakis E, Dunaif A. Insulin resistance and the polycystic ovary syndrome revisited: an update on mechanisms and implications. *Endocr Rev*. 2012;33(6):981-1030.
11. Rubin KH, Glintborg D, Nybo M, Abrahamsen B, Andersen M. Development and Risk Factors of Type 2 Diabetes in a Nationwide Population of Women With Polycystic Ovary Syndrome. *J Clin Endocrinol Metab*. 2017;102(10):3848-57.
12. Dunaif A. Perspectives in Polycystic Ovary Syndrome: From Hair to Eternity. *J Clin Endocrinol Metab*. 2016;101(3):759-68.
13. Zawadzki JKD, A. Diagnostic criteria for polycystic ovary syndrome; towards a rational approach. *Polycystic Ovary Syndrome*. Boston, Massachusetts: Blackwell Scientific; 1992. p. 377-84.
14. Rotterdam EA-SPcwg. Revised 2003 consensus on diagnostic criteria and long-term health risks related to polycystic ovary syndrome (PCOS). *Hum Reprod*. 2004;19(1):41-7.
15. Rotterdam EA-SPCWG. Revised 2003 consensus on diagnostic criteria and long-term health risks related to polycystic ovary syndrome. *Fertil Steril*. 2004;81(1):19-25.
16. Cimino I, Casoni F, Liu X, Messina A, Parkash J, Jamin SP, et al. Novel role for anti-Mullerian hormone in the regulation of GnRH neuron excitability and hormone secretion. *Nat Commun*. 2016;7:10055.

17. Tata B, Mimouni NEH, Barbotin AL, Malone SA, Loyens A, Pigny P, et al. Elevated prenatal anti-Mullerian hormone reprograms the fetus and induces polycystic ovary syndrome in adulthood. *Nat Med.* 2018;24(6):834-46.
18. Hayes MG, Urbanek M, Ehrmann DA, Armstrong LL, Lee JY, Sisk R, et al. Genome-wide association of polycystic ovary syndrome implicates alterations in gonadotropin secretion in European ancestry populations. *Nat Commun.* 2015;6:7502.
19. Chen ZJ, Zhao H, He L, Shi Y, Qin Y, Shi Y, et al. Genome-wide association study identifies susceptibility loci for polycystic ovary syndrome on chromosome 2p16.3, 2p21 and 9q33.3. *Nat Genet.* 2011;43(1):55-9.
20. Shi Y, Zhao H, Shi Y, Cao Y, Yang D, Li Z, et al. Genome-wide association study identifies eight new risk loci for polycystic ovary syndrome. *Nat Genet.* 2012;44(9):1020-5.
21. Day F, Karaderi T, Jones MR, Meun C, He C, Drong A, et al. Large-scale genome-wide meta-analysis of polycystic ovary syndrome suggests shared genetic architecture for different diagnosis criteria. *PLoS Genet.* 2018;14(12):e1007813.
22. Day FR, Hinds DA, Tung JY, Stolk L, Styrkarsdottir U, Saxena R, et al. Causal mechanisms and balancing selection inferred from genetic associations with polycystic ovary syndrome. *Nat Commun.* 2015;6:8464.
23. Castaldi PJ, Dy J, Ross J, Chang Y, Washko GR, Curran-Everett D, et al. Cluster analysis in the COPDGene study identifies subtypes of smokers with distinct patterns of airway disease and emphysema. *Thorax.* 2014;69(5):415-22.
24. Li L, Cheng WY, Glicksberg BS, Gottesman O, Tamler R, Chen R, et al. Identification of type 2 diabetes subgroups through topological analysis of patient similarity. *Sci Transl Med.* 2015;7(311):311ra174.
25. Tzeng CR, Chang YC, Chang YC, Wang CW, Chen CH, Hsu MI. Cluster analysis of cardiovascular and metabolic risk factors in women of reproductive age. *Fertil Steril.* 2014;101(5):1404-10.
26. Dewailly D, Alebic MS, Duhamel A, Stojanovic N. Using cluster analysis to identify a homogeneous subpopulation of women with polycystic ovarian morphology in a population of non-hyperandrogenic women with regular menstrual cycles. *Hum Reprod.* 2014;29(11):2536-43.
27. Daan NM, Koster MP, de Wilde MA, Dalmeijer GW, Evelein AM, Fauser BC, et al. Biomarker Profiles in Women with PCOS and PCOS Offspring; A Pilot Study. *PLoS One.* 2016;11(11):e0165033.
28. Huang CC, Tien YJ, Chen MJ, Chen CH, Ho HN, Yang YS. Symptom patterns and phenotypic subgrouping of women with polycystic ovary syndrome: association between endocrine characteristics and metabolic aberrations. *Hum Reprod.* 2015;30(4):937-46.
29. Dapas M, Sisk R, Legro RS, Urbanek M, Dunaif A, Hayes MG. Family-based quantitative trait meta-analysis implicates rare noncoding variants in DENND1A in polycystic ovary syndrome. *J Clin Endocrinol Metab.* 2019.

30. Legro RS, Driscoll D, Strauss JF, 3rd, Fox J, Dunaif A. Evidence for a genetic basis for hyperandrogenemia in polycystic ovary syndrome. *Proc Natl Acad Sci U S A*. 1998;95(25):14956-60.
31. McCarty CA, Chisholm RL, Chute CG, Kullo IJ, Jarvik GP, Larson EB, et al. The eMERGE Network: a consortium of biorepositories linked to electronic medical records data for conducting genomic studies. *BMC Med Genomics*. 2011;4:13.
32. Murtagh F, Legendre P. Ward's Hierarchical Agglomerative Clustering Method: Which Algorithms Implement Ward's Criterion? *J Classif*. 2014;31(3):274-95.
33. Strauss T, von Maltitz MJ. Generalising Ward's Method for Use with Manhattan Distances. *PLoS One*. 2017;12(1):e0168288.
34. Hennig C. Cluster-wise assessment of cluster stability. *Comput Stat Data An*. 2007;52(1):258-71.
35. Voight BF, Kang HM, Ding J, Palmer CD, Sidore C, Chines PS, et al. The metabochip, a custom genotyping array for genetic studies of metabolic, cardiovascular, and anthropometric traits. *PLoS Genet*. 2012;8(8):e1002793.
36. Price AL, Patterson NJ, Plenge RM, Weinblatt ME, Shadick NA, Reich D. Principal components analysis corrects for stratification in genome-wide association studies. *Nat Genet*. 2006;38(8):904-9.
37. O'Connell J, Gurdasani D, Delaneau O, Pirastu N, Ulivi S, Cocca M, et al. A general approach for haplotype phasing across the full spectrum of relatedness. *PLoS Genet*. 2014;10(4):e1004234.
38. Genomes Project C, Auton A, Brooks LD, Durbin RM, Garrison EP, Kang HM, et al. A global reference for human genetic variation. *Nature*. 2015;526(7571):68-74.
39. Das S, Forer L, Schonherr S, Sidore C, Locke AE, Kwong A, et al. Next-generation genotype imputation service and methods. *Nat Genet*. 2016;48(10):1284-7.
40. Marchini J, Howie B, Myers S, McVean G, Donnelly P. A new multipoint method for genome-wide association studies by imputation of genotypes. *Nat Genet*. 2007;39(7):906-13.
41. Willer CJ, Li Y, Abecasis GR. METAL: fast and efficient meta-analysis of genomewide association scans. *Bioinformatics*. 2010;26(17):2190-1.
42. Yang D, Jang I, Choi J, Kim MS, Lee AJ, Kim H, et al. 3DIV: A 3D-genome Interaction Viewer and database. *Nucleic Acids Res*. 2018;46(D1):D52-D7.
43. Shin H, Shi Y, Dai C, Tjong H, Gong K, Alber F, et al. TopDom: an efficient and deterministic method for identifying topological domains in genomes. *Nucleic Acids Res*. 2016;44(7):e70.
44. Venables WN, Ripley BD, Venables WN. *Modern applied statistics with S*. 4th ed. New York: Springer; 2002. xi, 495 p. p.
45. Tuomi T, Santoro N, Caprio S, Cai M, Weng J, Groop L. The many faces of diabetes: a disease with increasing heterogeneity. *Lancet*. 2014;383(9922):1084-94.
46. Skyler JS, Bakris GL, Bonifacio E, Darsow T, Eckel RH, Groop L, et al. Differentiation of diabetes by pathophysiology, natural history, and prognosis. *Diabetes*. 2017;66(2):241-55.

47. Ahlqvist E, Storm P, Käräjämäki A, Martinell M, Dorkhan M, Carlsson A, et al. Novel subgroups of adult-onset diabetes and their association with outcomes: a data-driven cluster analysis of six variables. *Lancet Diabetes Endocrinol.* 2018.
48. Udler MS, Kim J, von Grotthuss M, Bonàs-Guarch S, Cole JB, Chiou J, et al. Type 2 diabetes genetic loci informed by multi-trait associations point to disease mechanisms and subtypes: A soft clustering analysis. *PLOS Medicine.* 2018;15(9):e1002654.
49. Di Zazzo E, De Rosa C, Abbondanza C, Moncharmont B. PRDM Proteins: Molecular Mechanisms in Signal Transduction and Transcriptional Regulation. *Biology (Basel).* 2013;2(1):107-41.
50. Consortium GT. The Genotype-Tissue Expression (GTEx) project. *Nat Genet.* 2013;45(6):580-5.
51. Consortium F, the RP, Clst, Forrest AR, Kawaji H, Rehli M, et al. A promoter-level mammalian expression atlas. *Nature.* 2014;507(7493):462-70.
52. Carling T, Kim KC, Yang XH, Gu J, Zhang XK, Huang S. A histone methyltransferase is required for maximal response to female sex hormones. *Mol Cell Biol.* 2004;24(16):7032-42.
53. Liu L, Shao G, Steele-Perkins G, Huang S. The retinoblastoma interacting zinc finger gene RIZ produces a PR domain-lacking product through an internal promoter. *J Biol Chem.* 1997;272(5):2984-91.
54. Andreu-Vieyra C, Chen R, Matzuk MM. Conditional deletion of the retinoblastoma (Rb) gene in ovarian granulosa cells leads to premature ovarian failure. *Mol Endocrinol.* 2008;22(9):2141-61.
55. Yang QE, Nagaoka SI, Gwest I, Hunt PA, Oatley JM. Inactivation of Retinoblastoma Protein (Rb1) in the Oocyte: Evidence That Dysregulated Follicle Growth Drives Ovarian Teratoma Formation in Mice. *PLoS Genet.* 2015;11(7):e1005355.
56. Reader KL, Haydon LJ, Littlejohn RP, Juengel JL, McNatty KP. Booroola BMPR1B mutation alters early follicular development and oocyte ultrastructure in sheep. *Reprod Fertil Dev.* 2012;24(2):353-61.
57. Shimasaki S, Moore RK, Otsuka F, Erickson GF. The bone morphogenetic protein system in mammalian reproduction. *Endocr Rev.* 2004;25(1):72-101.
58. Estienne A, Pierre A, di Clemente N, Picard JY, Jarrier P, Mansanet C, et al. Anti-Mullerian hormone regulation by the bone morphogenetic proteins in the sheep ovary: deciphering a direct regulatory pathway. *Endocrinology.* 2015;156(1):301-13.
59. Yi SE, LaPolt PS, Yoon BS, Chen JY, Lu JK, Lyons KM. The type I BMP receptor Bmpr1B is essential for female reproductive function. *Proc Natl Acad Sci U S A.* 2001;98(14):7994-9.
60. Sugiura K, Su YQ, Eppig JJ. Does bone morphogenetic protein 6 (BMP6) affect female fertility in the mouse? *Biol Reprod.* 2010;83(6):997-1004.
61. Gorsic LK, Dapas M, Legro RS, Hayes MG, Urbanek M. Functional Genetic Variation in the Anti-Mullerian Hormone Pathway in Women with Polycystic Ovary Syndrome. *J Clin Endocrinol Metab.* 2019.

62. Gorsic LK, Kosova G, Werstein B, Sisk R, Legro RS, Hayes MG, et al. Pathogenic Anti-Mullerian Hormone Variants in Polycystic Ovary Syndrome. *J Clin Endocrinol Metab.* 2017;102(8):2862-72.
63. Pruitt KD, Tatusova T, Maglott DR. NCBI Reference Sequence (RefSeq): a curated non-redundant sequence database of genomes, transcripts and proteins. *Nucleic Acids Res.* 2005;33(Database issue):D501-4.
64. Shi W, Wymore RS, Wang HS, Pan Z, Cohen IS, McKinnon D, et al. Identification of two nervous system-specific members of the erg potassium channel gene family. *J Neurosci.* 1997;17(24):9423-32.
65. Hardy AB, Fox JE, Giglou PR, Wijesekara N, Bhattacharjee A, Sultan S, et al. Characterization of Erg K⁺ channels in alpha- and beta-cells of mouse and human islets. *J Biol Chem.* 2009;284(44):30441-52.
66. Muhlbauer E, Bazwinsky I, Wolgast S, Klemenz A, Peschke E. Circadian changes of ether-a-go-go-related-gene (Erg) potassium channel transcripts in the rat pancreas and beta-cell. *Cell Mol Life Sci.* 2007;64(6):768-80.
67. Wang D, Chu M, Wang F, Zhou A, Ruan M, Chen Y. A Genetic Variant in FIGN Gene Reduces the Risk of Congenital Heart Disease in Han Chinese Populations. *Pediatr Cardiol.* 2017;38(6):1169-74.
68. Wang D, Wang F, Shi KH, Tao H, Li Y, Zhao R, et al. Lower Circulating Folate Induced by a Fidgetin Intronic Variant Is Associated With Reduced Congenital Heart Disease Susceptibility. *Circulation.* 2017;135(18):1733-48.
69. Desbuquois B, Carre N, Burnol AF. Regulation of insulin and type 1 insulin-like growth factor signaling and action by the Grb10/14 and SH2B1/B2 adaptor proteins. *FEBS J.* 2013;280(3):794-816.
70. Kasus-Jacobi A, Perdereau D, Auzan C, Clauser E, Van Obberghen E, Mauvais-Jarvis F, et al. Identification of the rat adapter Grb14 as an inhibitor of insulin actions. *J Biol Chem.* 1998;273(40):26026-35.
71. Zhao W, Rasheed A, Tikkanen E, Lee JJ, Butterworth AS, Howson JMM, et al. Identification of new susceptibility loci for type 2 diabetes and shared etiological pathways with coronary heart disease. *Nat Genet.* 2017;49(10):1450-7.
72. Brodie A, Azaria JR, Ofran Y. How far from the SNP may the causative genes be? *Nucleic Acids Res.* 2016;44(13):6046-54.
73. Maher MC, Uricchio LH, Torgerson DG, Hernandez RD. Population genetics of rare variants and complex diseases. *Hum Hered.* 2012;74(3-4):118-28.
74. Park JH, Gail MH, Weinberg CR, Carroll RJ, Chung CC, Wang Z, et al. Distribution of allele frequencies and effect sizes and their interrelationships for common genetic susceptibility variants. *Proc Natl Acad Sci U S A.* 2011;108(44):18026-31.
75. Goring HH, Terwilliger JD, Blangero J. Large upward bias in estimation of locus-specific effects from genomewide scans. *Am J Hum Genet.* 2001;69(6):1357-69.
76. Kraft P. Curses--winner's and otherwise--in genetic epidemiology. *Epidemiology.* 2008;19(5):649-51; discussion 57-8.

77. McAllister JM, Modi B, Miller BA, Biegler J, Bruggeman R, Legro RS, et al. Overexpression of a DENND1A isoform produces a polycystic ovary syndrome theca phenotype. *Proc Natl Acad Sci U S A*. 2014;111(15):E1519-27.
78. Tee MK, Speek M, Legeza B, Modi B, Teves ME, McAllister JM, et al. Alternative splicing of DENND1A, a PCOS candidate gene, generates variant 2. *Mol Cell Endocrinol*. 2016;434:25-35.
79. Moran L, Teede H. Metabolic features of the reproductive phenotypes of polycystic ovary syndrome. *Hum Reprod Update*. 2009;15(4):477-88.
80. Fauser BC, Tarlatzis BC, Rebar RW, Legro RS, Balen AH, Lobo R, et al. Consensus on women's health aspects of polycystic ovary syndrome (PCOS): the Amsterdam ESHRE/ASRM-Sponsored 3rd PCOS Consensus Workshop Group. *Fertil Steril*. 2012;97(1):28-38 e25.
81. Carmina E, Koyama T, Chang L, Stanczyk FZ, Lobo RA. Does ethnicity influence the prevalence of adrenal hyperandrogenism and insulin resistance in polycystic ovary syndrome? *Am J Obstet Gynecol*. 1992;167(6):1807-12.
82. Guo M, Chen ZJ, Eijkemans MJ, Goverde AJ, Fauser BC, Macklon NS. Comparison of the phenotype of Chinese versus Dutch Caucasian women presenting with polycystic ovary syndrome and oligo/amenorrhoea. *Hum Reprod*. 2012;27(5):1481-8.
83. Louwers YV, Lao O, Fauser BC, Kayser M, Laven JS. The impact of self-reported ethnicity versus genetic ancestry on phenotypic characteristics of polycystic ovary syndrome (PCOS). *J Clin Endocrinol Metab*. 2014;99(10):E2107-16.
84. Essah PA, Nestler JE, Carmina E. Differences in dyslipidemia between American and Italian women with polycystic ovary syndrome. *J Endocrinol Invest*. 2008;31(1):35-41.
85. Dunaif A, Sorbara L, Delson R, Green G. Ethnicity and polycystic ovary syndrome are associated with independent and additive decreases in insulin action in Caribbean-Hispanic women. *Diabetes*. 1993;42(10):1462-8.
86. Engmann L, Jin S, Sun F, Legro RS, Polotsky AJ, Hansen KR, et al. Racial and ethnic differences in the polycystic ovary syndrome metabolic phenotype. *Am J Obstet Gynecol*. 2017;216(5):493 e1- e13.
87. Zhang W, Collins A, Gibson J, Tapper WJ, Hunt S, Deloukas P, et al. Impact of population structure, effective bottleneck time, and allele frequency on linkage disequilibrium maps. *Proc Natl Acad Sci U S A*. 2004;101(52):18075-80.
88. McCarthy MI. The importance of global studies of the genetics of type 2 diabetes. *Diabetes Metab J*. 2011;35(2):91-100.