

## Loci on chromosome 12q13.2 encompassing *ERBB3*, *PA2G4* and *RAB5B* are associated with polycystic ovary syndrome

Short Title: Loci on chromosome 12q13.2 and PCOS

Topics: polycystic ovary syndrome, theca cells, family cohort, single nucleotide polymorphisms, candidate genes

R. Alan Harris<sup>1</sup> [rharris1@bcm.edu](mailto:rharris1@bcm.edu)

Kellie J. Archer<sup>2</sup> [archer.43@osu.edu](mailto:archer.43@osu.edu)

Mark O. Goodarzi<sup>3</sup> [mark.goodarzi@cshs.org](mailto:mark.goodarzi@cshs.org)

Timothy P. York<sup>4,5</sup> [tpyork@vcu.edu](mailto:tpyork@vcu.edu)

Jeffrey Rogers<sup>1</sup> [jr13@bcm.edu](mailto:jr13@bcm.edu)

Andrea Dunaif<sup>6</sup> [andrea.dunaif@mssm.edu](mailto:andrea.dunaif@mssm.edu)

Jan McAllister<sup>7</sup> [jxm63@psu.edu](mailto:jxm63@psu.edu)

Jerome F Strauss III<sup>5, 8, \*</sup> [jerome.strauss@pennmedicine.upenn.edu](mailto:jerome.strauss@pennmedicine.upenn.edu)

<sup>1</sup> Human Genome Sequencing Center and Department of Molecular and Human Genetics, Baylor College of Medicine, Houston, TX 77030 USA

<sup>2</sup> Division of Biostatistics, College of Public Health, The Ohio State University, Columbus, OH 43210 USA

<sup>3</sup> Division of Endocrinology, Diabetes, and Metabolism, Cedars-Sinai Medical Center, Los Angeles, CA 90048 USA

<sup>4</sup> Department of Human and Molecular Genetics, Virginia Commonwealth University, Richmond, VA 23298 USA

<sup>5</sup> Department of Obstetrics and Gynecology, Virginia Commonwealth University School of Medicine, Richmond, Virginia 23298 USA

<sup>6</sup> Division of Endocrinology, Diabetes and Bone Disease, Department of Medicine, Icahn School of Medicine at Mount Sinai, New York, NY 10029 USA

<sup>7</sup> Department of Pathology, Penn State Hershey College of Medicine, Hershey, PA 17033 USA

<sup>8</sup> Department of Obstetrics and Gynecology, University of Pennsylvania Perelman School of Medicine, Philadelphia, Pennsylvania, 19104 USA

\* Corresponding author

1 **Abstract**

2 Polycystic ovary syndrome (PCOS) is characterized by hyperandrogenemia of ovarian theca cell  
3 origin. Here we report the significant association of 15 single nucleotide polymorphisms (SNPs),  
4 identified by whole exome sequencing (WES). DNA was isolated from well-characterized theca  
5 cell preparations from women of European ancestry with PCOS (N=9) and elevated androgen  
6 production in vitro and from normal ovulatory women (N=7). Of the SNPs, 10 are located within  
7 150 kb on chromosome 12q13.2. This region contains three plausible PCOS candidate genes  
8 (*ERBB3/PA2G4/RAB5B*), two of which (*ERBB3* and *RAB5B*) have been identified in GWAS of  
9 PCOS. None of the SNPs individually had a significant association with PCOS or in vivo androgen  
10 levels when evaluated in an independent cohort (n=318) of families with one or more  
11 daughters with PCOS, but a haplotype consisting of the minor alleles of three of the SNPS  
12 (rs773121, rs773123 and rs812826) was found preferentially in women with PCOS and elevated  
13 androgen levels (p=0.0583). Moreover, the three minor alleles in this haplotype were  
14 significantly associated with anti-Mullerian Hormone (AMH) levels, a marker of follicular  
15 reserve and follicular maturation. Two of the three SNP minor alleles are predicted to have  
16 significant functional consequences (rs773123 a missense SNP in *ERBB3*, and rs812826, a SNP in  
17 the *PA2G4* promoter). Notably, *PA2G4* encodes a protein that interacts with the *ERBB3*  
18 cytoplasmic domain, which is also the domain where the missense variant resides. These  
19 findings provide support for the contribution and probable functional significance of loci on  
20 chromosome 12q13.2 to the pathophysiology underlying PCOS.

21

22

## 23 **Author Summary**

24 Polycystic ovary syndrome (PCOS) is the most common endocrine disorder of women of  
25 reproductive age. We identified 15 single nucleotide polymorphisms (SNPs) associated with  
26 androgen production in theca cells from normal ovulatory women and women with PCOS. Of  
27 these SNPs, 10 are within a 150 kbp region of chromosome 12 including 9 that form a  
28 haplotype. This region contains three PCOS candidate genes (*ERBB3/PA2G4/RAB5B*). These  
29 SNPs were further examined in an independent cohort of families with one or more daughters  
30 with PCOS. A haplotype consisting of the minor alleles of three of the SNPS was found  
31 preferentially in women with PCOS. Furthermore, the three minor alleles in this haplotype were  
32 significantly associated with anti-Mullerian Hormone (AMH) levels, a marker of follicular  
33 reserve and maturation. Two SNPs in the chromosome 12 haplotype were likely to have  
34 functional consequences based on genomic context, suggesting that they affect *ERBB3* and  
35 *PA2G4* interactions or *PA2G4* expression.

## 36 **Introduction**

37 Polycystic ovary syndrome (PCOS) is the most common endocrine disorder of women of  
38 reproductive age. PCOS is characterized by anovulatory infertility, hyperandrogenemia and  
39 metabolic disturbance [1]. PCOS is a complex genetic disorder and approximately 20  
40 susceptibility loci have been reproducibly associated in genome wide association studies  
41 (GWAS) [2]. However, the functional significance of many of the genes in these loci with  
42 respect to PCOS phenotypes is largely unknown [2].

43 We examined the potential role of genes implicated in PCOS GWAS in controlling  
44 androgen production by ovarian theca cells from normal cycling women and women with PCOS,

45 since hyperandrogenemia of ovarian origin is implicated in the pathophysiology of ovarian  
46 dysfunction and metabolic disturbances in PCOS [3]. Here we describe SNPs and a haplotype  
47 located in a 150 kb region of chromosome 12q13.2 that contains plausible PCOS genes  
48 (*ERBB3/PA2G4/RAB5B*), and variants that could have a causal role in promoting PCOS  
49 phenotypes.

## 50 **Results**

### 51 **Whole Exome Sequencing Reveals SNPs on Chromosome 12 Associated with Thecal Cell** 52 **Androgen Production.**

53 Table 1 lists the 18 PCOS candidate genes [4–8] interrogated in the whole exome  
54 sequencing study. These genes were selected from published GWAS and meta-analyses,  
55 representing genes in loci associated with PCOS in women of European ancestry. The variants  
56 identified by WES in genes of interest are presented in Supplemental Table S1. The gene locus,  
57 chromosome position, rs number, transcript GenBank accession number, the reference and  
58 minor alleles, location in the gene, and number of minor alleles detected are provided.

59 PCOS theca cells produced significantly more DHEA, the major androgen secreted by  
60 theca cell cultures, compared to theca cells from normal ovulatory women under basal  
61 conditions, and when challenged with forskolin, which activates adenylate cyclase and mimics  
62 the action of luteinizing hormone (LH) (Table 2).

63 Table 3 presents the number of variants analyzed in each gene of interest after filtering  
64 to remove duplicate entries and variants that were homogeneous across all cell preparations.  
65 An allele-based Wilcoxon rank sum test was then performed on each of the filtered SNPs  
66 (N=252) to detect different levels of forskolin-stimulated DHEA production by variant (Table 4).

67 Among the 15 variants that were found to be significant at a  $P$ -value  $< 0.05$ , 10 were located on  
68 chromosome 12. This over representation of chromosome 12 was highly significant (Fisher's  
69 exact test,  $P=0.000002$ ).

70 The variants on chromosome 12 showing significant effects include 10 SNPs, 6 in the  
71 *ERBB3* gene (rs7297125, rs12817471, rs2229046, rs773123, rs812826, rs773121), and 4  
72 (rs67594137, rs11171713, rs11550558, rs7963590) in *RAB5B/SUOX*. The latter genes overlap  
73 each other on opposite DNA strands. Notably, the same PCOS theca cell preparations (MC01,  
74 MC03, MC27 representing one third of the PCOS sample) contained the minor allele of these  
75 SNPs in a heterozygous state, suggesting linkage disequilibrium (LD) and a large effect size. The  
76 LD was not unexpected since the variants are located within a 150 kb stretch of chromosome  
77 12q13.2. Moreover, using LDlink programs to investigate linkage disequilibrium in European  
78 populations, we found that one SNP, rs7297175, is independent of the other 9, while 6 SNPs  
79 (rs67594137, rs11171713, rs11550558, rs7963590, rs12817471, rs2229046) formed one LD  
80 group ( $r^2$  0.77 to 0.97 among the SNPs) and 3 SNPs (rs773123, rs812826, rs773121) formed  
81 another LD group ( $r^2$  1.0) (Figure 1). We then used LDlink to identify the haplotypes present in  
82 the two haplotype blocks, finding two haplotypes in each block, one consisting of the major  
83 allele of the constituent SNPs and the other containing the minor alleles (Figure 2).

#### 84 **Analysis of Whole Genome Sequences from a Cohort of European Ancestry Women.**

85 We examined the chromosome 12 SNPs in whole genome sequencing data from a  
86 cohort of 318 individuals of European ancestry from 77 families with one or more daughters  
87 with PCOS [9]. PLINK [10] transmission disequilibrium tests (TDT) on the individual SNPs did not  
88 identify any of them as significantly associated with PCOS/HA (hyperandrogenemia) affection

89 status [11] (Table S2). However, rs773123, rs812826, and rs773121, corresponding to one of  
90 the haplotype groups identified in the initial samples, had p values ranging from 0.1011 to  
91 0.1404 compared to p values ranging from 0.4652 to 0.8575 for the other SNPs. Based on the  
92 lower range of p values for rs773123, rs812826, and rs773121, we examined the possibility they  
93 formed a haplotype with association to PCOS/HA using the Family-Based Association Tests  
94 (FBAT) [12] HBAT test which is the haplotype version of the association test. Based on the HBAT  
95 test, the haplotype consisting of the minor allele for rs773123, rs812826, and rs773121 and the  
96 major allele for the remaining SNPs was found preferentially in women with PCOS and elevated  
97 androgen levels ( $p = 0.0583$ ) (Table S3).

98 Association analyses of levels of total testosterone (T), DHEAS (dehydroepiandrosterone  
99 sulfate), sex hormone binding globulin (SHBG), luteinizing hormone (LH), follicle stimulating  
100 hormone (FSH), and anti-Mullerian hormone (AMH) were performed using the PLINK family-  
101 based association test for quantitative traits (QFAM). After 100,000 permutations for empirical  
102 p-value correction of QFAM total results, the only associations significant at  $p < 0.01$  were AMH  
103 associations with rs773123 ( $P_{emp} = 0.00105$ ), rs812826 ( $P_{emp} = 0.0009$ ), and rs773121 ( $P_{emp} =$   
104  $0.00242$ ). The minor alleles of these same SNPs were present in the haplotype that approached  
105 significance in the affection status FBAT HBAT test. The quantitative trait data for T, DHEAS,  
106 SHBG, LH and FSH were reported as part of our previous study using family-based WGS analyses  
107 [9]. The AMH median (25th to 75th percentiles) value (ng/mL) for affected women was 10.02  
108 (6.64-14.31) and for unaffected women, it was 3.02 (1.06-4.77). No quantitative analyses of  
109 hormone levels using FBAT HBAT were significant.

## 110 **Functional Annotations of SNPs**

111 To examine the potential functional impact of the chromosome 12 SNPs, we annotated  
112 them with CADD PHRED scores [13]. A CADD PHRED score of 10 predicts a SNP is among the  
113 10% most functional changes in the human genome while a CADD PHRED score of 20 indicates  
114 the change is among the 1% most functional changes. Three of the SNPs have a CADD PHRED  
115 score greater than 10 including two in the HBAT identified haplotype (Table S4). rs773123  
116 (CADD PHRED 24.6) is a missense SNP in *ERBB3* at amino acid position 1119 which is in the  
117 cytoplasmic topological domain  
118 ([https://www.uniprot.org/uniprot/P21860#subcellular\\_location](https://www.uniprot.org/uniprot/P21860#subcellular_location)). rs773121 (CADD PHRED 13.2)  
119 is in the promoter (Ensembl regulatory feature ENSR00000052477) for *PA2G4* which is a  
120 transcriptional co-repressor of androgen receptor-regulated genes [14]. *PA2G4* interacts with  
121 the cytoplasmic domain of *ERBB3* [15] which contains the rs773123 missense SNP. The *PA2G4*  
122 Ensembl regulatory feature ENSR00000052477 is equivalent to the GeneHancer [16]  
123 GH12J056102 regulatory element. This GeneHancer regulatory element is classified as having  
124 both promoter and enhancer functions and two of the genes for which it has a predicted  
125 enhancer function are *ERBB3* and *RAB5B*. These interactions can be viewed in the UCSC  
126 Genome Browser “Interactions between GeneHancer regulatory elements and genes” track  
127 (<https://genome.ucsc.edu/s/Rharris1/chr12.PCOS>). Since these two SNPs occur in a haplotype  
128 with a p value approaching significance, there could be some interplay between SNP rs773121’s  
129 regulatory effect on the expression of *PA2G4* and *ERBB3* and SNP rs773123’s effect on the  
130 interaction between the *ERBB3* cytoplasmic domain and *PA2G4*. Additionally, *PA2G4* and *ERBB3*  
131 are targets of transcription factor, *ZNF217*, another PCOS candidate gene identified by GWAS

132 ([https://maayanlab.cloud/Harmonizome/gene\\_set/ZNF217/ENCODE+Transcription+Fac](https://maayanlab.cloud/Harmonizome/gene_set/ZNF217/ENCODE+Transcription+Fac)  
133 [tor+Targets](#)) [17].

## 134 Discussion

135 We found evidence of an association between androgen (DHEA) production by cultures  
136 of human theca cells, a reflection of theca cell endocrine activity, and fifteen SNPs, including  
137 ten in a haplotype on chromosome 12. This haplotype was found preferentially in women with  
138 PCOS and elevated androgen levels ( $p = 0.0583$ ) in an independent family cohort. Two minor  
139 alleles in the chromosome 12 haplotype were likely to have functional consequences based on  
140 CADD scores and genomic context, suggesting that they affect *ERBB3* and *PA2G4* interactions  
141 (*rs773123*) or *PA2G4* expression (*rs773121*). Moreover, a recent study [18] utilizing a different  
142 methodologic approach, colocalization analysis, identified the same region on chromosome 12  
143 encompassing *ERBB3*, *PA2G4*, *RAB5B* and *SUOX*, as containing potential PCOS disease-  
144 mediating genes.

145 *ERBB3* is a component of the EGF family of signal transduction factors. It forms  
146 heterodimers with other ERBB family members, including *ERBB2*, which is the receptor for  
147 neuregulins, which are produced by both theca cells and granulosa cells in response to LH [19].  
148 Neuregulin-1 (NRG-1) also plays a role in the proliferation of Leydig cells, an androgen  
149 producing cell type in the male gonad that is functionally analogous to theca cells in the female  
150 gonad. Thus, the haplotype identified in this report encompasses genes that are implicated in  
151 the regulation/function of androgen producing cells. Moreover, both *ERBB3* and *RAB5B* have  
152 been associated with metabolic phenotypes that are related to PCOS, including glucose  
153 metabolism/insulin resistance [4,5,20].



154 GWAS conducted on Han Chinese identified *RAB5B/SUOX* as PCOS candidate loci [7].  
155 The minor alleles of the SNPs we identified have very different allele frequencies in non-Finnish  
156 European and East Asian populations with the exception of rs11550558, being low frequency in  
157 East Asians (<0.05%) and 7-11% in Europeans (Table S5). A recent case-control study of SNPs  
158 encoding the 3'-UTR of the *RAB5B* gene conducted in Han Chinese revealed highly significant  
159 associations with PCOS phenotypes with large effect sizes [21]. These findings support our  
160 conclusions that 12q13.2 contains important genetic determinants of PCOS, and that the  
161 specific SNPs that influence thecal androgen production are population-specific. The SNPs may  
162 impact the level of expression of the candidate genes, accounting for variation in cell/tissue  
163 function, including the endocrine functions of thecal cells and granulosa cells. For example,  
164 *ERBB3* is expressed by granulosa cells [19], which produce AMH, a peptide factor that  
165 influences growth of follicles [22]. Circulating AMH levels are elevated in women with PCOS,  
166 suggesting that the chromosome 12 haplotype plays a role in regulating both theca cell  
167 (androgen) and granulosa (AMH) cell function [23].

168 *RAB5B* is involved in intracellular trafficking of endosomes including those derived from  
169 the plasma membrane [24]. We have previously shown that *RAB5B* is co-localized in  
170 compartments containing *DENND1A.V2*, another PCOS GWAS candidate gene associated with  
171 hyperandrogenemia [25,26]. *DENND1A.V2* has an N-terminal guanine nucleotide exchange  
172 function and a clathrin-binding domain, putting it at the nexus with plasma membrane signaling  
173 proteins like the *ERBB* family of proteins. *DENND1A.V2* is translocated into the nucleus of PCOS  
174 theca cells along with *RAB5B*, suggesting a role in regulation of expression of genes involved in  
175 androgen synthesis [26]. *SUOX*, which overlaps the *RAB5B* gene, encodes a mitochondrial

176 sulfite oxidase, which has not been linked to PCOS, thecal cell function or steroidogenesis.

177 Thus, SNPs in this gene are likely to be irrelevant to the PCOS phenotypes.

178 In summary, our findings provide support for functional contributions of genes on  
179 chromosome 12q13.2, including *ERBB3*, *PA2G4*, and *RAB5B*, to ovarian cell dysfunction in PCOS.  
180 Our findings align with in silico colocalization analyses implicating these same genes in PCOS  
181 pathogenesis.

## 182 **Materials and Methods**

### 183 Theca cell preparations and culture

184 Human theca interna tissue was obtained from follicles of women undergoing  
185 hysterectomy, following informed consent under a protocol approved by the Institutional  
186 Review Board of The Pennsylvania State University College of Medicine. As a standard of care,  
187 oophorectomies were performed during the luteal phase of the cycle. Theca cells from normal  
188 cycling and PCOS follicles were isolated and grown as we have as previously reported in  
189 detail[27–29]. PCOS and normal ovarian tissue came from age-matched women, 38–41 years  
190 old. The diagnosis of PCOS was made according to National Institutes of Health (NIH) consensus  
191 guidelines [30,31] which include hyperandrogenemia/hyperandrogenism and oligo-ovulation  
192 and the exclusion of other causes of hyperandrogenemia (e.g. 21-hydroxylase deficiency,  
193 Cushing’s syndrome, and adrenal or ovarian tumors). All of the PCOS theca cell preparations  
194 studied came from ovaries of women with fewer than six menses per year and elevated serum  
195 total testosterone or bioavailable testosterone levels [32–35]. Each of the PCOS ovaries  
196 contained multiple subcortical follicles of less than 10 mm in diameter. The control (normal)  
197 theca cell preparations came from ovaries of fertile women with normal menstrual histories,

198 menstrual cycles of 21–35 days, and no clinical signs of hyperandrogenism. Neither PCOS nor  
199 normal subjects were receiving hormonal medications at the time of surgery. Indications for  
200 surgery were dysfunctional uterine bleeding, endometrial cancer, and pelvic pain. Experiments  
201 comparing PCOS and normal theca were performed using fourth-passage (31–38 population  
202 doublings) theca cells isolated from individual size-matched follicles obtained from age-  
203 matched subjects, in the absence of in vivo stimulation. The use of fourth-passage cells allowed  
204 us to perform multiple experiments from the same patient population, and were propagated  
205 from frozen stocks of second passage cells in the media described above. The passage  
206 conditions and split ratios for all normal and PCOS cells were identical. These studies were  
207 approved by the Human Subjects Protection Offices of Virginia Commonwealth University) and  
208 Penn State College of Medicine.

209         Nine cell preparations obtained from women with PCOS (MC01, MC21, MC09\_B, MC03,  
210 MC10, MC16, MC26, MC27, MC190) and seven from normal ovulating women (MC62\_B, MC02,  
211 MC06, MC31, MC38, MC40, MC50) were studied. All subjects were unrelated and of European  
212 ancestry. The cells were characterized by their production of dehydroepiandrosterone (DHEA),  
213 the major androgen synthesized by these cells, under basal conditions or stimulated with  
214 forskolin (20  $\mu$ M) for 16 h. DHEA was quantified by ELISA assays (DRG, Springfield, NJ) and  
215 production (pmol) was normalized to cell number ( $10^6$  cells) determined at the end of the  
216 culture period.

#### 217 Whole Exome Sequence Analysis of Normal and PCOS Theca Cells DNA

218         The DNA samples were subjected to whole exome sequencing at 100 millions reads  
219 providing 100 $\times$  coverage using the Agilent SureSelect 51M capture kit with Illumina HiSeq 2000

220 sequencing, in conjunction with BGI Americas. Raw sequence data for each individual were  
221 mapped to the human reference genome (build GRCh37/hg19) using the BWA-MEM algorithm  
222 of Burrows-Wheeler Aligner (v 0.7.12) (H. Li, 2013). This was followed by a series of pre-  
223 processing steps—marking duplicates, realignment around indels and base quality recalibration.  
224 PCR duplicates were marked within the aligned reads using Picard tools.  
225 (<http://picard.sourceforge.net>) Next, mapping artifacts around indels were cleaned up using  
226 the RealignerTargetCreator, the IndelRealigner and the LeftAlignIndels walkers of the Genome  
227 Analysis ToolKit (GATK) [37,38]. Inaccurate / biased base quality scores were recalibrated using  
228 the BaseRecalibrator, the AnalyzeCovariates and the PrintReads walkers of GATK, which use  
229 machine learning to model these errors empirically and adjust the quality scores accordingly.

### 230 Linkage disequilibrium

231 LDlink (<https://analysistools.cancer.gov/LDlink/?tab=home>) was used to identify  
232 haplotypes in Europeans. LDlink accesses data from 1000 Genomes in a suite of tools that  
233 allows determination of linkage disequilibrium (LD) and haplotypes. We used the LDlink SNPclip  
234 tool to examine LD among the 9 SNPs associated with DHEA response, using an  $r^2$  cutoff of 0.5.  
235 This identified two LD groups, which were then used to identify haplotypes using LDhap.

### 236 Statistical analysis

237 The Wilcoxon rank sum test was used to compare age, basal DHEA, and forskolin-  
238 stimulated DHEA production between the PCOS and control samples. The WES data were  
239 subjected to the following filters. We retained genetic variants having a unique combination of  
240 Gene ID, Chromosome, Position, Variant ID, Reference and Alternate Allele. Additionally,  
241 variants that were homogeneous across all samples (i.e., no sample displayed the minor allele

242 or all samples displayed the minor allele (N=21)) were removed, leaving 441 variants for  
243 statistical analysis.

244 For each variant, a Wilcoxon rank sum test was used to compare those with and without  
245 the variant with respect to forskolin-stimulated DHEA production. In order to apply a statistical  
246 comparison, a minimum of two samples per group were required so that the set of variants was  
247 restricted to 252 variants. P values of <0.05 were considered significant. We did not apply  
248 Bonferroni correction because we are testing a restricted set of genetic variants in robust loci  
249 for PCOS.

#### 250 SNP analyses of a PCOS cohort

251 We analyzed the 10 SNPs on chromosome 12 that were significantly associated with  
252 thecal cell androgen production in whole genome sequencing data from a family based PCOS  
253 cohort [9]. The study was approved by the Institutional Review Boards of Northwestern  
254 University Feinberg School of Medicine, Penn State Health Milton S. Hershey Medical Center,  
255 and Brigham and Women's Hospital. Written informed consent was obtained from all subjects  
256 prior to the study. The cohort consisted of 318 individuals of European ancestry from 77  
257 families with one or more daughters with PCOS. Among the index cases and sisters (n=171), the  
258 following phenotypes were identified: PCOS (T>58 ng/dl and/or uT>15 ng/dl and  $\leq 8$   
259 menses/year) (n=90); Hyperandrogenemic (HA) (T>58 ng/dl and/or uT>15 ng/dl and regular  
260 menses (every 27-35 days)) (n=5); Unaffected (n=76). The women were ages 14 to 49 years.  
261 Women were assigned affected status if they fulfilled criteria for PCOS or HA, as we have done  
262 in our previous family-based genetic analyses [11]. All anthropometric and hormonal data,  
263 except for circulating AMH levels, have been previously reported [9].

264 Sequencing of the cohort was performed using the Complete Genomics, Inc. platform.  
265 Sequence reads were aligned to the human reference genome (GRCh37/hg19) and variants  
266 were called using the CGI AssemblyPipeline version 2.0. The SNPs were analyzed individually  
267 using the PLINK v1.90 [10] transmission disequilibrium test (TDT) based on PCOS affection  
268 status. An individual was considered affected if they had a phenotype of PCOS or  
269 Hyperandrogenic. The PLINK family-based association test for quantitative traits (QFAM) was  
270 used to examine associations of quantitative traits with the SNPs measured in the index cases  
271 and sisters. The sample size was T (n=162), DHEAS (n=161), SHBG (n=160), LH (n=162) and FSH  
272 (n=162). There was insufficient sample for AMH assays in some subjects, thus, the sample size  
273 with AMH values was smaller (n=59). QFAM total results were empirically corrected using  
274 100,000 permutations. The haplotypes containing the SNPs were analyzed using the Family-  
275 Based Association Tests (FBAT) v2.0.3 [12] HBAT function which is the haplotype version of the  
276 association test. FBAT HBAT was performed using both the PCOS affection status and  
277 quantitative traits. Potential functional consequences of the SNPs were examined using  
278 Combined Annotation-Dependent Depletion (CADD) v1.6 [13] and GeneHancer [16].  
279

280 **References**

- 281
- 282 1. Diamanti-Kandarakis E, Dunaif A. Insulin resistance and the polycystic ovary syndrome  
283 revisited: an update on mechanisms and implications. *Endocr Rev.* 2012;33: 981–1030.  
284 doi:10.1210/ER.2011-1034
  - 285 2. Dapas M, Dunaif A. Deconstructing a Syndrome: Genomic Insights into PCOS Causal  
286 Mechanisms and Classification. *Endocr Rev.* 2022 [cited 19 Apr 2022].  
287 doi:10.1210/ENDREV/BNAC001
  - 288 3. Legro RS, Driscoll D, Strauss JF, Fox J, Dunaif A. Evidence for a genetic basis for  
289 hyperandrogenemia in polycystic ovary syndrome. *Proc Natl Acad Sci U S A.* 1998;95:  
290 14956–14960. doi:10.1073/PNAS.95.25.14956
  - 291 4. Day FR, Hinds DA, Tung JY, Stolk L, Styrkarsdottir U, Saxena R, et al. Causal mechanisms  
292 and balancing selection inferred from genetic associations with polycystic ovary  
293 syndrome. *Nat Commun.* 2015;6. doi:10.1038/NCOMMS9464
  - 294 5. Day F, Karaderi T, Jones MR, Meun C, He C, Drong A, et al. Large-scale genome-wide  
295 meta-analysis of polycystic ovary syndrome suggests shared genetic architecture for  
296 different diagnosis criteria. *PLoS Genet.* 2018;14. doi:10.1371/JOURNAL.PGEN.1007813
  - 297 6. Hayes MG, Urbanek M, Ehrmann DA, Armstrong LL, Lee JY, Sisk R, et al. Genome-wide  
298 association of polycystic ovary syndrome implicates alterations in gonadotropin secretion  
299 in European ancestry populations. *Nat Commun.* 2015/08/19. 2015;6: 7502.  
300 doi:10.1038/ncomms8502
  - 301 7. Shi Y, Zhao H, Shi Y, Cao Y, Yang D, Li Z, et al. Genome-wide association study identifies  
302 eight new risk loci for polycystic ovary syndrome. *Nat Genet.* 2012;44: 1020–1025.  
303 doi:10.1038/NG.2384
  - 304 8. Chen ZJ, Zhao H, He L, Shi Y, Qin Y, Shi Y, et al. Genome-wide association study identifies  
305 susceptibility loci for polycystic ovary syndrome on chromosome 2p16.3, 2p21 and  
306 9q33.3. *Nat Genet.* 2011;43: 55–59. doi:10.1038/NG.732
  - 307 9. Dapas M, Sisk R, Legro RS, Urbanek M, Dunaif A, Hayes MG. Family-based quantitative  
308 trait meta-analysis implicates rare noncoding variants in DENND1A in polycystic ovary  
309 syndrome. *J Clin Endocrinol Metab.* 2019;104: 3835–3850. doi:10.1210/JC.2018-02496
  - 310 10. Purcell S, Neale B, Todd-Brown K, Thomas L, Ferreira MAR, Bender D, et al. PLINK: a tool  
311 set for whole-genome association and population-based linkage analyses. *Am J Hum*  
312 *Genet.* 2007;81: 559–575. doi:10.1086/519795
  - 313 11. Urbanek M, Legro RS, Driscoll DA, Azziz R, Ehrmann DA, Norman RJ, et al. Thirty-seven  
314 candidate genes for polycystic ovary syndrome: strongest evidence for linkage is with  
315 follistatin. *Proc Natl Acad Sci U S A.* 1999;96: 8573–8578. doi:10.1073/PNAS.96.15.8573
  - 316 12. Horvath S, Xu X, Laird NM. The family based association test method: strategies for  
317 studying general genotype--phenotype associations. *Eur J Hum Genet.* 2001;9: 301–306.  
318 doi:10.1038/SJ.EJHG.5200625
  - 319 13. Rentzsch P, Witten D, Cooper GM, Shendure J, Kircher M. CADD: predicting the  
320 deleteriousness of variants throughout the human genome. *Nucleic Acids Res.* 2019;47:  
321 D886–D894. doi:10.1093/NAR/GKY1016

- 322 14. Zhang Y, Akinmade D, Hamburger AW. The ErbB3 binding protein Ebp1 interacts with  
323 Sin3A to repress E2F1 and AR-mediated transcription. *Nucleic Acids Res.* 2005;33: 6024–  
324 6033. doi:10.1093/NAR/GKI903
- 325 15. Yoo JY, Wang XW, Rishi AK, Lessor T, Xia XM, Gustafson TA, et al. Interaction of the  
326 PA2G4 (EBP1) protein with ErbB-3 and regulation of this binding by heregulin. *Br J*  
327 *Cancer.* 2000;82: 683–690. doi:10.1054/BJOC.1999.0981
- 328 16. Fishilevich S, Nudel R, Rappaport N, Hadar R, Plaschkes I, Stein TI, et al. GeneHancer:  
329 genome-wide integration of enhancers and target genes in GeneCards. Database  
330 (Oxford). 2017;2017. doi:10.1093/DATABASE/BAX028
- 331 17. Krig SR, Miller JK, Frieze S, Beckett LA, Neve RM, Farnham PJ, et al. ZNF217, a candidate  
332 breast cancer oncogene amplified at 20q13, regulates expression of the ErbB3 receptor  
333 tyrosine kinase in breast cancer cells. *Oncogene.* 2010;29: 5500–5510.  
334 doi:10.1038/ONC.2010.289
- 335 18. Censin JC, Bovijn J, Holmes M v, Lindgren CM. Colocalization analysis of polycystic ovary  
336 syndrome to identify potential disease-mediating genes and proteins. *Eur J Hum Genet.*  
337 2021/03/06. 2021;29: 1446–1454. doi:10.1038/s41431-021-00835-8
- 338 19. Chowdhury I, Branch A, Mehrabi S, Ford BD, Thompson WE. Gonadotropin-Dependent  
339 Neuregulin-1 Signaling Regulates Female Rat Ovarian Granulosa Cell Survival.  
340 *Endocrinology.* 2017;158: 3647–3660. doi:10.1210/EN.2017-00065
- 341 20. Jones MR, Brower MA, Xu N, Cui J, Mengesha E, Chen YDI, et al. Systems Genetics  
342 Reveals the Functional Context of PCOS Loci and Identifies Genetic and Molecular  
343 Mechanisms of Disease Heterogeneity. *PLoS Genet.* 2015;11.  
344 doi:10.1371/JOURNAL.PGEN.1005455
- 345 21. Yu J, Ding C, Guan S, Wang C. Association of single nucleotide polymorphisms in the  
346 RAB5B gene 3'UTR region with polycystic ovary syndrome in Chinese Han women. *Biosci*  
347 *Rep.* 2019/05/01. 2019;39. doi:10.1042/BSR20190292
- 348 22. Dumont A, Robin G, Catteau-Jonard S, Dewailly D. Role of Anti-Müllerian Hormone in  
349 pathophysiology, diagnosis and treatment of Polycystic Ovary Syndrome: a review.  
350 *Reprod Biol Endocrinol.* 2015;13. doi:10.1186/S12958-015-0134-9
- 351 23. Gorsic LK, Kosova G, Werstein B, Sisk R, Legro RS, Hayes MG, et al. Pathogenic Anti-  
352 Müllerian Hormone Variants in Polycystic Ovary Syndrome. *J Clin Endocrinol Metab.*  
353 2017;102: 2862–2872. doi:10.1210/JC.2017-00612
- 354 24. Gulappa T, Clouser CL, Menon KM. The role of Rab5a GTPase in endocytosis and post-  
355 endocytic trafficking of the hCG-human luteinizing hormone receptor complex. *Cell Mol*  
356 *Life Sci.* 2011;68: 2785–2795. doi:10.1007/s00018-010-0594-1
- 357 25. McAllister JM, Legro RS, Modi BP, Strauss 3rd JF. Functional genomics of PCOS: from  
358 GWAS to molecular mechanisms. *Trends Endocrinol Metab.* 2015/01/21. 2015;26: 118–  
359 124. doi:10.1016/j.tem.2014.12.004
- 360 26. Kulkarni R, Teves ME, Han AX, McAllister JM, Strauss 3rd JF. Colocalization of Polycystic  
361 Ovary Syndrome Candidate Gene Products in Theca Cells Suggests Novel Signaling  
362 Pathways. *J Endocr Soc.* 2019/11/15. 2019;3: 2204–2223. doi:10.1210/js.2019-00169
- 363 27. Wickenheisser JK, Nelson-Degrave VL, McAllister JM. Dysregulation of cytochrome P450  
364 17alpha-hydroxylase messenger ribonucleic acid stability in theca cells isolated from  
365 women with polycystic ovary syndrome. *J Clin Endocrinol Metab.* 2005;90: 1720–1727.



- 366 Available:  
367 [http://www.ncbi.nlm.nih.gov/entrez/query.fcgi?cmd=Retrieve&db=PubMed&dopt=Citation&list\\_uids=15598676](http://www.ncbi.nlm.nih.gov/entrez/query.fcgi?cmd=Retrieve&db=PubMed&dopt=Citation&list_uids=15598676)  
368
- 369 28. Wickenheisser JK, Biegler JM, Nelson-Degrave VL, Legro RS, Strauss 3rd JF, McAllister JM.  
370 Cholesterol side-chain cleavage gene expression in theca cells: augmented transcriptional  
371 regulation and mRNA stability in polycystic ovary syndrome. *PLoS One*. 2012/11/17.  
372 2012;7: e48963. doi:10.1371/journal.pone.0048963
- 373 29. Nelson-Degrave VL, Wickenheisser JK, Hendricks KL, Asano T, Fujishiro M, Legro RS, et al.  
374 Alterations in mitogen-activated protein kinase kinase and extracellular regulated kinase  
375 signaling in theca cells contribute to excessive androgen production in polycystic ovary  
376 syndrome. *Mol Endocrinol*. 2005;19: 379–390. Available:  
377 [http://www.ncbi.nlm.nih.gov/entrez/query.fcgi?cmd=Retrieve&db=PubMed&dopt=Citation&list\\_uids=15514033](http://www.ncbi.nlm.nih.gov/entrez/query.fcgi?cmd=Retrieve&db=PubMed&dopt=Citation&list_uids=15514033)  
378
- 379 30. Legro RS, Arslanian SA, Ehrmann DA, Hoeger KM, Murad MH, Pasquali R, et al. Diagnosis  
380 and treatment of polycystic ovary syndrome: an endocrine society clinical practice  
381 guideline. *J Clin Endocrinol Metab*. 2013/10/24. 2013;98: 4565–4592.  
382 doi:10.1210/jc.2013-2350
- 383 31. Azziz R, Carmina E, Chen Z, Dunaif A, Laven JS, Legro RS, et al. Polycystic ovary syndrome.  
384 *Nat Rev Dis Primers*. 2016/08/12. 2016;2: 16057. doi:10.1038/nrdp.2016.57
- 385 32. Nelson-DeGrave VL, Wickenheisser JK, Cockrell JE, Wood JR, Legro RS, Strauss 3rd JF, et  
386 al. Valproate potentiates androgen biosynthesis in human ovarian theca cells.  
387 *Endocrinology*. 2004;145: 799–808. Available:  
388 [http://www.ncbi.nlm.nih.gov/entrez/query.fcgi?cmd=Retrieve&db=PubMed&dopt=Citation&list\\_uids=14576182](http://www.ncbi.nlm.nih.gov/entrez/query.fcgi?cmd=Retrieve&db=PubMed&dopt=Citation&list_uids=14576182)  
389
- 390 33. Nelson VL, Legro RS, Strauss 3rd JF, McAllister JM. Augmented androgen production is a  
391 stable steroidogenic phenotype of propagated theca cells from polycystic ovaries. *Mol*  
392 *Endocrinol*. 1999;13: 946–957. Available:  
393 [http://www.ncbi.nlm.nih.gov/entrez/query.fcgi?cmd=Retrieve&db=PubMed&dopt=Citation&list\\_uids=10379893](http://www.ncbi.nlm.nih.gov/entrez/query.fcgi?cmd=Retrieve&db=PubMed&dopt=Citation&list_uids=10379893)  
394
- 395 34. Wickenheisser JK, Nelson-DeGrave VL, Quinn PG, McAllister JM. Increased cytochrome  
396 P450 17alpha-hydroxylase promoter function in theca cells isolated from patients with  
397 polycystic ovary syndrome involves nuclear factor-1. *Mol Endocrinol*. 2004;18: 588–605.  
398 Available:  
399 [http://www.ncbi.nlm.nih.gov/entrez/query.fcgi?cmd=Retrieve&db=PubMed&dopt=Citation&list\\_uids=14684846](http://www.ncbi.nlm.nih.gov/entrez/query.fcgi?cmd=Retrieve&db=PubMed&dopt=Citation&list_uids=14684846)  
400
- 401 35. Wickenheisser JK, Quinn PG, Nelson VL, Legro RS, Strauss 3rd JF, McAllister JM.  
402 Differential activity of the cytochrome P450 17alpha-hydroxylase and steroidogenic  
403 acute regulatory protein gene promoters in normal and polycystic ovary syndrome theca  
404 cells. *J Clin Endocrinol Metab*. 2000;85: 2304–2311. Available:  
405 [http://www.ncbi.nlm.nih.gov/entrez/query.fcgi?cmd=Retrieve&db=PubMed&dopt=Citation&list\\_uids=10852468](http://www.ncbi.nlm.nih.gov/entrez/query.fcgi?cmd=Retrieve&db=PubMed&dopt=Citation&list_uids=10852468)  
406
- 407 36. Li H. Aligning sequence reads, clone sequences and assembly contigs with BWA-MEM.  
408 2013 [cited 19 Dec 2017]. Available: <http://arxiv.org/abs/1303.3997>

- 409 37. McKenna A, Hanna M, Banks E, Sivachenko A, Cibulskis K, Kernytsky A, et al. The Genome  
410 Analysis Toolkit: a MapReduce framework for analyzing next-generation DNA sequencing  
411 data. *Genome Res.* 2010;20: 1297–1303. doi:10.1101/GR.107524.110
- 412 38. Depristo MA, Banks E, Poplin R, Garimella K v., Maguire JR, Hartl C, et al. A framework for  
413 variation discovery and genotyping using next-generation DNA sequencing data. *Nat*  
414 *Genet.* 2011;43: 491–501. doi:10.1038/NG.806
- 415
- 416
- 417

<b>Locus</b>	<b>Population</b>	<b>GWAS Reference</b>
<i>C9orf3</i>	Han, European	Shi et al., 2012; Hayes et al., 2015
<i>DENND1A</i>	Han, European	Chen et al., 2011; Shi et al., 2012; Day et al., 2018
<i>ERBB2</i>	European	Day et al., 2018
<i>ERBB3</i>	European	Day et al., 2018
<i>ERBB4</i>	European	Day et al., 2015, 2018
<i>FSHB</i>	European	Day et al., 2015; Hayes et al., 2015
<i>GATA4/NEIL2</i>	European	Hayes et al., 2015
<i>KRR1</i>	European	Day et al., 2015, 2018
<i>MAPRE1</i>	European	Day et al., 2018
<i>PLGRKT</i>	European	Day et al., 2018
<i>RAB5B/SUOX</i>	Han, European	Shi et al., 2012; Day et al., 2018
<i>RAD50</i>	European	Day et al., 2015, 2018
<i>THADA</i>	Han, European	Chen et al., 2011; Shi et al., 2012; Day et al., 2015, 2018
<i>TOX3</i>	Han, European	Shi et al., 2012; Day et al., 2018
<i>YAP1</i>	Han, European	Shi et al., 2012; Day et al., 2015, 2018
<i>ZBTB16</i>	European	Day et al., 2018

418

419 **Table 1: PCOS Candidate Loci Interrogated.**

420 **A**

Sample	Age	Basal DHEA pmol/10 <sup>6</sup> cells	Forskolin-DHEA pmol/10 <sup>6</sup> cells
<b>Normal</b>			
MC02	41-45	40.37 ± 2.37	183.32 ± 2.57
MC06	36-40	45.23 ± 3.99	304.90 ± 30.22
MC31	36-40	24.50 ± 2.09	139.64 ± 15.94
MC38	36-40	52.87 ± 22.60	578.52 ± 5.34
MC40	41-45	46.95 ± 3.60	156.40 ± 14.57
MC50	36-40	35.84 ± 2.65	147.88 ± 8.11
MC62	36-40	31.95 ± 0.58	273.46 ± 35.13
<b>PCOS</b>			
MC01	36-40	386.71 ± 49.65	6504.54 ± 358.27
MC03	26-30	1281.66 ± 214.17	6412.01 ± 558.38
MC09	36-40	799.64 ± 29.87	5326.77 ± 328.37
MC10	31-35	196.43 ± 10.69	2005.88 ± 105.57
MC16	31-35	148.43 ± 9.76	2714.01 ± 100.84
MC21	31-35	814.00 ± 21.21	5350.00 ± 212.13
MC26	26-30	439.96 ± 66.17	1928.96 ± 120.18
MC27	31-35	837.22 ± 78.87	5880.64 ± 1131.21
MC190	41-45	128.60 ± 11.25	3749.16 ± 365.41

421

422 **B**

	Normal	PCOS	p
N	7	9	
Age	37 (36, 41)	34 (30, 41)	0.101
Basal DHEA production	40.37 (24.5, 52.87)	439.96 (128.6, 1281.66)	<0.001
Forskolin-stimulated DHEA production	183.32 (139.64, 578.52)	5326.77 (1928.96, 6504.54)	<0.001

423

424

425 **Table 2: Androgen production by PCOS and normal thecal cells employed in this study. (A)**

426 DHEA production by normal and PCOS theca cell preparations employed in this study and  
 427 summary statistics. Cells were cultured for 16 h with or without forskolin (20 μM) and DHEA  
 428 production was assessed by immunoassay normalized to cell number. Values presented are  
 429 means (S.D.) from triplicate cultures for each preparation. **(B)** Median and range (minimum,  
 430 maximum) with P-value from Wilcoxon rank sum test.

431

432

<b>Gene</b>	<b>Variants</b>
C9orf3	31
DENND1A	48
ERBB2	14
ERBB3	23
ERBB4	65
FSHB	4
GATA4	42
KRR1	17
MAPRE1	5
NEIL2	15
PLGRKT	13
RAB5B	11
RAD50	24
SUOX	4
THADA	74
TOX3	21
YAP1	13
ZBTB16	17

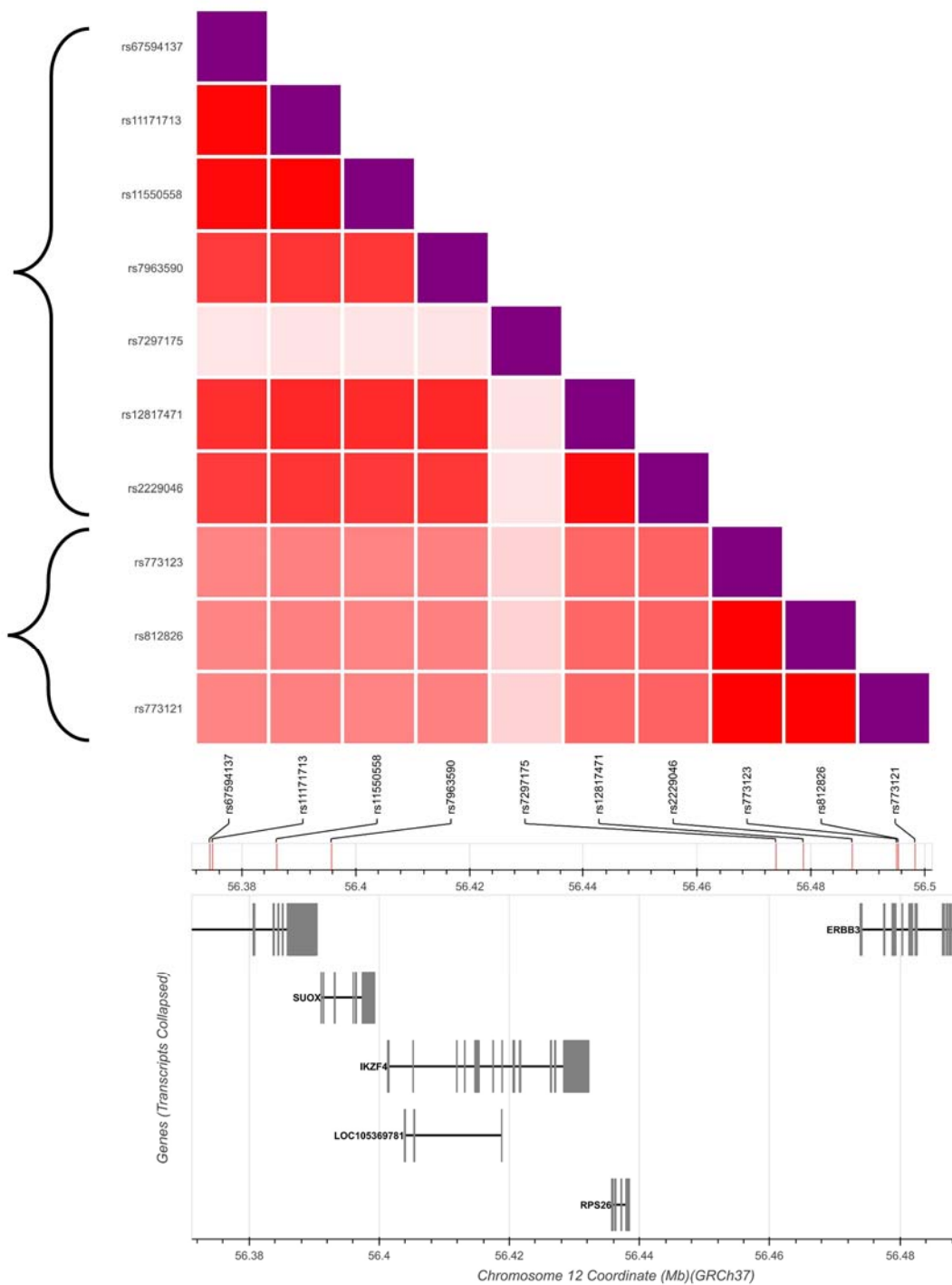
433

434 **Table 3. Number of variants analyzed by gene after filtering steps applied**

Gene	Chromosome	Position	dbSNP ID	Major Allele	Minor Allele	P-value
THADA	chr2	43519977	rs35720761	C	T	0.013278
THADA	chr2	43736171	rs6544669	T	C	0.007692
ERBB4	chr2	212530466	rs6725181	C	T	0.029670
ERBB4	chr2	212587321	rs13002712	G	A	0.033333
DENND1A	chr9	126304695	rs748994	T	G	0.039286
RAB5B	chr12	56374318	rs67594137	C	T	0.003571
RAB5B	chr12	56374803	rs11171713	G	A	0.003571
RAB5B	chr12	56386076	rs11550558	A	G	0.003571
SUOX	chr12	56395689	rs7963590	G	A	0.003571
ERBB3	chr12	56473808	rs7297175	T	C	0.041783
ERBB3	chr12	56478607	rs12817471	G	A	0.003571
ERBB3	chr12	56487201	rs2229046	T	C	0.003571
ERBB3	chr12	56494998	rs773123	A	T	0.003571
ERBB3	chr12	56495306	rs812826	C	T	0.003571
ERBB3	chr12	56498241	rs773121	G	A	0.003571

435

436 **Table 4. Variants having significantly different forskolin-stimulated DHEA when comparing**  
437 **samples with the minor allele to those with only the major allele.** SNPs significantly associated  
438 with thecal androgen production. Nominal p values are presented.



439  
440

441 **Fig 1. Linkage disequilibrium among chromosome 12 SNPs in Individuals of European**  
 442 **Ancestry** The plot displays linkage disequilibrium as  $r^2$  between each pair of SNPs. The brackets  
 443 outline the two linkage disequilibrium groups ( $r^2 > 0.5$ ). Note that rs7297175 is independent of  
 444 the other 9 SNPs and is not included in either haplotype block. Darker red indicates higher  
 445 linkage disequilibrium. SNP locations and genes in the region are displayed at bottom.

RS Number	Position (GRCh37)	Allele Frequencies	Haplotypes	
rs67594137	chr12:56374318	C=0.926, T=0.074	C	T
rs11171713	chr12:56374803	G=0.928, A=0.072	G	A
rs11550558	chr12:56386076	A=0.927, G=0.073	A	G
rs7963590	chr12:56395689	G=0.93, A=0.07	G	A
rs12817471	chr12:56478607	G=0.925, A=0.075	G	A
rs2229046	chr12:56487201	T=0.927, C=0.073	T	C
<b>Haplotype Count</b>			<b>922</b>	<b>59</b>
<b>Haplotype Frequency</b>			<b>0.9165</b>	<b>0.0586</b>

RS Number	Position (GRCh37)	Allele Frequencies	Haplotypes	
rs773123	chr12:56494998	A=0.888, T=0.112	A	T
rs812826	chr12:56495306	C=0.888, T=0.112	C	T
rs773121	chr12:56498241	G=0.888, A=0.112	G	A
<b>Haplotype Count</b>			<b>893</b>	<b>113</b>
<b>Haplotype Frequency</b>			<b>0.8877</b>	<b>0.1123</b>

446  
447

448 **Fig 2. Haplotypes in the chromosome 12 region of interest.** The first haplotype block consists  
449 of 6 SNPs and the second consists of 3 SNPs. Each block contains a common haplotype and a  
450 rare haplotype.

451



452 **Supporting information captions**

453

454 **Supplemental Table 1.xlsx**

455

456 **Table S1. WES results for selected PCOS candidate genes.** The table presents the gene name,  
457 chromosome assignment, nucleotide position of the detected variant, rs number, major and  
458 minor alleles, location and/or variant effect on coding sequence, transcript, and distribution of  
459 minor alleles among the different theca cell preparations designated by their MC number (see  
460 below), with number of homozygous (left), heterozygous (middle) followed by the total number  
461 of minor alleles (right) detected.

462

463 **Supplemental.docx**

464

465 **Table S2.** PLINK transmission disequilibrium tests (TDT) based on affection status for the  
466 individual chromosome 12 SNPs.

467

468 **Table S3.** Family-Based Association Tests (FBAT) HBAT test based on affection status for  
469 chromosome 12 SNPs.

470

471 **Table S4.** CADD PHRED scores and a subset of annotation details for chromosome 12 SNPs.

472

473 **Table S5. Minor Allele Frequencies for the 12q13.2 SNPs.** Minor allele frequencies for non-  
474 Finnish Europeans, East Asians and African Americans were extracted from GnomAD  
475 (<https://gnomad.broadinstitute.org>). rs7297175 is independent of the other 9 SNPs and is not  
476 included in either haplotype block.

477