

1 A genome-wide association study of polycystic ovary syndrome identified from
2 electronic health records

3 Yanfei Zhang^{1,*,**}, MD, PhD; Kevin Ho^{2,#,*}, MD; Jacob M. Keaton^{3,4}, PhD; Dustin N. Hartzel⁵, BS;
4 Felix Day⁶, PhD; Anne E. Justice⁷, PhD; Navya S. Josyula⁷, MS; Sarah A. Pendergrass^{7, #}, PhD;
5 Ky'Era Actkins^{4,8,9}, BS; Lea K. Davis^{4,8,10,11}, PhD; Digna R. Velez Edwards^{4,11,12}, PhD; Brody
6 Holohan¹³, PhD; Andrea Ramirez¹⁴, MD, MS; Ian B. Stanaway¹⁵, PhD; David R. Crosslin¹⁵, PhD;
7 Gail P. Jarvik¹⁶, MD, PhD; Patrick Sleiman¹⁷, PhD; Hakon Hakonarson¹⁷, MD, PhD; Marc S.
8 Williams¹, MD; Ming Ta Michael Lee^{1,**}, PhD.

9

10 1. Genomic Medicine Institute, Geisinger, Danville, PA, USA

11 2. Kidney Research Institute, Geisinger, Danville, PA, USA

12 3. Division of Epidemiology, Department of Medicine; Institute for Medicine and Public Health;
13 Vanderbilt University Medical Center, Nashville, TN, USA

14 4. Vanderbilt Genetics Institute, Vanderbilt University, Nashville, TN, USA

15 5. Phenomic Analytics and Clinical Data Core, Geisinger, PA, USA

16 6. The International PCOS Consortium

17 7. Department of Population Health Sciences, Geisinger, Danville, PA, USA

18 8. Division of Genetic Medicine, Department of Medicine, Vanderbilt University Medical Center,
19 Nashville, TN, USA

20 9. Department of Microbiology, Immunology, and Physiology, Meharry Medical College,
21 Nashville, TN, USA

22 10. Department of Psychiatry and Behavioral Sciences; Vanderbilt University Medical Center,
23 Nashville, TN, USA

24 11. Department of Biomedical Informatics, Data Sciences Institute, Vanderbilt University
25 Medical Center, Nashville, TN, USA

26 12. Division of Quantitative Science, Department of Obstetrics and Gynecology, Vanderbilt
27 University Medical Center, Nashville, TN, USA

28 13. Marshfield Clinic Research Institute, Marshfield, WI, USA

29 14. Department of Medicine, Vanderbilt University Medical Center, Nashville, TN, USA

30 15. R.

31 16. Departments of Medicine (Medical Genetics) and Genome Sciences, School of Medicine,
32 University of Washington, Seattle, WA, USA

33 17. Children's Hospital of Philadelphia, Philadelphia, PA, USA

34

35 * The first two authors contribute equally.

36 # Both Kevin Ho (currently employed by Sanofi Genzyme) and Sarah A. Pendergrass (currently
37 employed by Genentech) worked on this study while employed by Geisinger.

38 ** Co-Corresponding authors:

39 Yanfei Zhang, y Zhang@geisinger.edu; Ming Ta Michael Lee, mlee2@geisinger.edu;

40 Genomic Medicine Institute; Weis Center for Research, Geisinger, 100 N. Academy Ave.

41 Danville, PA 17822, USA

42 Abstract

43 **Background:** Polycystic ovary syndrome (PCOS) is the most common endocrine disorder
44 affecting women of reproductive age. Previous studies have identified genetic variants
45 associated with PCOS identified by different diagnostic criteria. The Rotterdam Criteria is the
46 broadest and able to identify the most PCOS cases.

47 **Objectives:** To identify novel associated genetic variants, we extracted PCOS cases and controls
48 from the electronic health records (EHR) based on the Rotterdam Criteria and performed a
49 genome-wide association study (GWAS).

50 **Study Design:** We developed a PCOS phenotyping algorithm based on the Rotterdam criteria
51 and applied it to three EHR-linked biobanks to identify cases and controls for genetic study. In
52 discovery phase, we performed individual GWAS using the Geisinger's MyCode and the
53 eMERGE cohorts, which were then meta-analyzed. We attempted validation of the significantly
54 association loci ($P < 1 \times 10^{-6}$) in the BioVU cohort. All association analyses used logistic regression,
55 assuming an additive genetic model, and adjusted for principal components to control for
56 population stratification. An inverse-variance fixed effect model was adopted for meta-
57 analyses. Additionally, we examined the top variants to evaluate their associations with each
58 criterion in the phenotyping algorithm. We used STRING to identify protein-protein interaction
59 network.

60 **Results:** We identified 2,995 PCOS cases and 53,599 controls in total (2,742 cases and 51,438
61 controls from the discovery phase; 253 cases and 2,161 controls in the validation phase). GWAS
62 identified one novel genome-wide significant variant rs17186366 (OR=1.37 [1.23,1.54],
63 $P = 2.8 \times 10^{-8}$) located near *SOD2*. Additionally, two loci with suggestive association were also

64 identified: rs113168128 (OR=1.72 [1.42,2.10], $P=5.2 \times 10^{-8}$), an intronic variant of *ERBB4* that is
65 independent from the previously published variants, and rs144248326 (OR=2.13 [1.52,2.86],
66 $P=8.45 \times 10^{-7}$), a novel intronic variant in *WWTR1*. In the further association tests of the top 3
67 SNPs with each criterion in the PCOS algorithm, we found that rs17186366 was associated with
68 polycystic and hyperandrogenism, while rs11316812 and rs144248326 were mainly associated
69 with oligomenorrhea or infertility. Besides *ERBB4*, we also validated the association with
70 *DENND1A1*.

71 **Conclusion:** Through a discovery-validation GWAS on PCOS cases and controls identified from
72 EHR using an algorithm based on Rotterdam criteria, we identified and validated a novel
73 association with variants within *ERBB4*. We also identified novel associations nearby *SOD2* and
74 *WWTR1*. These results suggest the eGFR and Hippo pathways in the disease etiology. With
75 previously identified PCOS-associated loci *YAP1*, the *ERBB4-YAP1-WWTR1* network implicates
76 the epidermal growth factor receptor and the Hippo pathway in the multifactorial etiology of
77 PCOS.

78

79 [Keywords](#)

80 EGFR pathway; Electronic Health Record; *ERBB4*; Hippo pathway; Polycystic Ovary Syndrome;

81 *SOD2*; *WWTR1*

82 Introduction

83 Polycystic ovary syndrome (PCOS) is the most common endocrine disorder that affects women
84 of reproductive age ¹. PCOS is characterized by the three main features: dysregulation of the
85 menstrual cycle; elevated levels of androgenic hormones; and multiple cysts of the ovaries from
86 which the name of the condition is derived. Other features include hirsutism in a “male”
87 pattern, acne, increased skin pigment sometimes associated with skin tags, and weight gain.
88 Three criteria to identify women with PCOS have been proposed: the National Institutes of
89 Health (NIH) Criteria, the Rotterdam Criteria, and the Androgen Excess and PCOS Society (AE-
90 PCOS) criteria. The NIH criteria requires both hyperandrogenism and oligomenorrhea ²; the
91 Rotterdam criteria requires at least two of the three phenotypes: hyperandrogenism, oligo-
92 ovulation, and polycystic ovaries ³; and the AE-PCOS requires both hyperandrogenism and
93 ovarian dysfunction ⁴. The Rotterdam criteria is more inclusive than the other two which
94 increases its sensitivity, thus, the prevalence of PCOS estimated by the Rotterdam criteria is 15-
95 20% compared to the 7-12% generated by the other two criteria ^{5, 6}.

96 Heritability estimated for PCOS ranges 38-71% by twin studies ⁷, with a polygenic genetic
97 architecture and complex inheritance pattern ^{8, 9}. Recent large-scale genome-wide association
98 studies (GWAS) have identified 19 loci associated with PCOS in women with European or East
99 Asian ancestries, including *ERBB4*, *YAP1* and *DENND1A* that were replicated in European and
100 Asian ancestral groups, providing additional evidence for the polygenic architecture of PCOS ¹⁰⁻
101 ¹⁶. These studies adopted different diagnosis criteria, including PCOS cases diagnosed based on
102 NIH or Rotterdam criteria, or self-reported information. Shared genetic architecture for PCOS

103 using the different diagnosis criteria or self-reported were also identified by genetic correlation
104 analyses^{13,16}.

105 The healthcare system-based biobanks with genetic data linked to the electronic health record
106 (EHR) data enable new opportunities for genomic discovery research¹⁷. Examples include
107 Geisinger's MyCode[®] Community Health Initiative (MyCode)¹⁸, BioVU at Vanderbilt University
108 ^{19, 20}, and the electronic Medical Records and Genomics (eMERGE) Network, a nationwide
109 consortium of multiple medical institutions that link DNA biobanks to EHRs²¹. These
110 multidimensional data are important resources for development of phenotype algorithms,
111 genetic discoveries, and clinical implementation²²⁻²⁴. Phenotyping algorithms to identify cases
112 and controls for various diseases from EHR have been developed²⁵. Such approaches are
113 critical for genetic studies as they integrate data from different EHR systems derived using the
114 same phenotype definition that has been rigorously evaluated to define the performance
115 characteristics in order to reduce case selection bias and heterogeneity among different
116 studies.

117 In this study, we aim to develop an EHR algorithm for PCOS based on the Rotterdam criteria to
118 identified PCOS cases and controls in multiple cohort and perform a GWAS to identify genetic
119 variants associated with PCOS.

120 Cohort and Methods

121 Cohorts

122 The discovery cohorts were identified from the Geisinger MyCode Community Health Initiative
123 (MyCode) Phase I ~ II and eMERGE Phase III. All MyCode participants provide written consent

124 allowing their clinical and genomic data to be used for health-related research^{18, 26}. The
125 eMERGE Phase III includes 83,717 individuals recruited from 12 study sites with demographics,
126 diagnosis information based on ICD codes, and genotyping data²⁴. The replication cohort was
127 selected from BioVU, Vanderbilt University's EHR-linked biorepository^{19, 20}. This study was
128 waived for a standard institutional review board (IRB) review based on the use of deidentified
129 EHR and genetic data from all sites. We received approval from the Geisinger MyCode
130 Governing Board, the eMERGE coordinating center and the BioVU Review Committee and IRB
131 to conduct this genetic study. Since both Geisinger and Vanderbilt are eMERGE sites,
132 participants in MyCode and BioVU who were included in the eMERGE data were excluded from
133 the site-specific analysis to avoid double-counting.

134 [PCOS EHR algorithm based on the Rotterdam Criteria](#)

135 **Figure 1** illustrates the sample selection and analytic strategy of this study. The Geisinger PCOS
136 EHR algorithm based on the Rotterdam diagnosis criteria was first developed to identify PCOS
137 cases and controls from the EHR data. The three criteria that were used in the algorithm to
138 represent different aspects of PCOS are: 1) Polycystic (C1): having diagnosis codes of polycystic
139 ovarian syndrome and/or polycystic ovaries; 2) Hyperandrogenic (C2): having diagnosis codes
140 for hyperandrogenism or hyperandrogenism-related clinical signs or hyperandrogenemia
141 determined by testosterone measurements; 3) Reproductive (C3): having diagnosis codes for
142 oligomenorrhea, amenorrhea, infertility and oligo- or anovulation. **Supplementary Table 1**
143 provides details and the inclusion and exclusion ICD codes and laboratory tests for each
144 criterion. PCOS cases were patients that met at least 2 of the 3 criteria with an index age
145 between 18 to 45. Controls were those who did not have any components of the three criteria,

146 and whose current age was older than the median age of the cases (38 years in this study) to
147 increase the specificity for the controls. This algorithm was then applied to the Geisinger and
148 eMERGE cohorts for the discovery GWAS.

149 [Discovery GWAS and meta-analyses](#)

150 MyCode Phase I and Phase II samples were genotyped and imputed to HRC.r1-1 EUR reference
151 genome (GRCh37 build) separately using the Michigan Imputation Server as previously
152 described ²⁷. Variants with imputation info score > 0.7 were included for analyses. eMERGE
153 samples were genotyped at each study site and imputed to HRC.r1-1 EUR reference genome in
154 multiple batches using the Michigan Imputation Server. Data were processed centrally and
155 harmonized as previously described ²⁴. Variants with average info score >0.3 were included.
156 Samples with a genotyping rate below 95% were excluded. SNPs with a <99% call rate, minor
157 allele frequency (MAF) of <1% and a significant deviation from the Hardy-Weinberg equilibrium
158 ($P < 1 \times 10^{-7}$) were removed from analyses. One of the paired individuals with first- or second-
159 degree relatedness were removed. Finally, there were 7,595,111 SNPs, 6,747,339 SNPs, and
160 5,648,769 SNPs from MyCode Phase I (1,141 cases and 18,788 controls), MyCode Phase II (594
161 cases and 9,024 controls), and eMERGE III (1,007 cases and 23,626 controls) included for GWAS.
162 For study-specific GWAS, we used fixed effects logistic regression, assuming an additive genetic
163 model, adjusted for index age and the first six principal components (PCs) to account for
164 population stratification for the MyCode Phase I and II cohorts; additionally, we adjusted for
165 the eMERGE III study sites. EasyQC ²⁸ was employed to harmonize the alleles and data format
166 for GWAS summary statistics from discovery studies before performing a fixed effect inverse

167 variance weighted meta-analysis using METAL²⁹. PLINK 1.9³⁰ was used to calculate PCs,
168 relatedness and to perform GWAS.

169 [Replication for the top variants](#)

170 Top associated variants with $P < 1 \times 10^{-6}$ from discovery meta-analysis were further evaluated in
171 an independent PCOS cohort identified based on the same algorithm from BioVU. We identified
172 253 cases and 2161 controls. Genotypes were generated using the Illumina Infinium Expanded
173 Multi-Ethnic Genotyping Array. The same imputation, quality control measures and association
174 protocols were applied for the replication study. We also queried the summary statistics of the
175 meta-analyses from the PCOS consortium (without the 23andMe data) for the associations of
176 these top variants¹⁶. Criteria for replication is $P < 0.05$, directionally consistent in the replication
177 GWAS, or $P < 5 \times 10^{-8}$ in the combined meta-analyses.

178 [Power calculation](#)

179 We evaluated the power for our study conservatively assuming a significance level of $P < 5 \times 10^{-8}$
180 for GWAS, and a PCOS prevalence of 8%. Given the current PCOS case number of 2995, we have
181 80% power to identify an associated variant with a MAF of 1% and an OR (odds ratio) > 2.01 ; or
182 a MAF at 2% and an OR > 1.68 ; or a MAF $> 8\%$ and an OR > 1.34 .

183 [Functional genomics exploration](#)

184 The Variant Effect Predictor was used for variant annotation³¹. The Functional Mapping and
185 Annotation was used in the default setting to generate independent loci and associated
186 pathways³². Open Targets Genetics is an online portal used in this study to query the
187 associated genes, phenome-wide association studies (PheWAS), and the expression

188 quantitative trait loci (eQTLs) of the top associated variants³³. The PheWAS data in Open Target
189 Genetics includes the results from UK Biobank GWAS and the GWAS catalog. STRING was used
190 to identify the protein-protein interaction network.

191 Results

192 Identification and characterization of PCOS cases and controls from EHR data

193 **Figure S1** illustrates the details of sample ascertainment using the Rotterdam-based algorithm
194 in the Geisinger and eMERGE cohorts. Only non-Hispanic whites were included in the MyCode
195 and BioVU samples; All races were included in the eMERGE samples, 75% were of European
196 American, 17% were of African American and 8% were other race/ethnicity (**Supplementary**
197 **table 2**). **Table 1** summarizes the numbers and characteristics of the identified cases and
198 controls. The proportion of patients with polycystic ovaries was around 40% of the PCOS cases
199 identified in the eMERGE and BioVU data; while this number is over 88% in the Geisinger
200 MyCode Phase I and II data. The Geisinger cohorts also have lower hyperandrogenic features
201 than eMERGE and BioVU cohorts. Over 90% of the patients had reproductive issues in all the
202 three cohorts. Cases showed higher BMI than controls in the MyCode and BioVU samples but
203 not in the eMERGE samples.

204 Discovery and replication of the genetic variants associated with the risk of PCOS

205 Twenty independent loci were identified with $P < 1 \times 10^{-5}$ in the discovery meta-analysis of the
206 MyCode and eMERGE cohorts (**Supplementary table 3**). Manhattan plots for the meta-analysis
207 and the three discovery studies are shown in **Figure 2A** and **Figure S2**. Variants with $P < 1 \times 10^{-6}$
208 were then examined in an independent cohort identified using the same algorithm from BioVU.

209 **Figure 2B** lists the association of the top three independent SNPs in the discovery and
210 replication cohorts. rs17186366, a novel association located in the promotor flanking region
211 near *SOD2* and *LOC101929142* reached genome-wide significance (OR=1.37 [1.23,1.54],
212 $P=2.8 \times 10^{-8}$) in the combined meta-analyses of discovery and replication. We also identified an
213 intronic variant of *ERBB4*, rs113168128, with near genome-wide significance (OR=1.72
214 [1.42,2.10], $P=5.2 \times 10^{-8}$). This SNP is independent from the previously reported *ERBB4* variant
215 rs2178575 ($r^2 = 0.001$)¹⁶. A low-frequency intronic variant of *WWTR1* rs144248326 (MAF =
216 1.01%), was identified and was also close to genome-wide significance level (OR=2.13
217 [1.52,2.86], $P=8.45 \times 10^{-7}$). The regional association plots for these three *loci* are shown in **Figure**
218 **S3**. We also examined the top associations for African Americans in the eMERGE datasets. Only
219 rs113168128 in *ERBB4* passed the standard quality check. This variant has higher MAF and a
220 slightly smaller OR with nominal significance (MAF=0.063, OR=1.64, $P=0.0106$).

221 We did not observe significant associations of these variants in the PCOS consortium meta-
222 analyses¹⁶. We also examined the associations of previously reported PCOS loci in our meta-
223 analyses result. Only variants in *DENND1A1* (rs9696009, rs10818854, rs10986105) were
224 replicated with the same direction and similar effect size ($P < 0.05$; **Supplementary table 4**).

225 The functional genomics exploration found rs17186366 near *SOD2* associated with menarche
226 ($P=6.6 \times 10^{-5}$), rs113168128 in *ERBB4* associated with depressed affect ($P=2.1 \times 10^{-5}$)³⁴
227 (**Supplementary table 5**). All of the top three SNPs were found to be associated with
228 phenotypes related to the nervous system, or to mental or behavioral disorders
229 (**Supplementary table 5**). None of these SNPs were found to be eQTLs in any tissue. The

230 protein-protein interaction network found both *ERBB4* and *WWTR1* interact with YAP1 (Figure
231 2C), which also associated with PCOS in both European and Han Chinese^{12, 13, 16}.

232 Association of the top three SNPs with each PCOS criterion

233 **Table 2** summarized the association results for the top three variants with each of the three
234 criteria in the PCOS algorithm that represent the different aspects of PCOS based on the
235 Rotterdam criteria. rs17186366 strongly associated with the polycystic and hyperandrogenic
236 traits, while the other two SNPs in *ERBB4* and *WWTR1* are mainly associated with the
237 reproductive trait as the more significant association P values and larger effect sizes are
238 observed for the variants and the corresponding traits (**Table 2**).

239 Discussion

240 Principal findings

241 In this study, we developed an EHR algorithm based on the Rotterdam criteria for PCOS and
242 identified cases and controls from three biobank cohorts. Through a discovery-validation
243 GWAS, we identified rs17186366 near *SOD2*, a novel genome-wide significant signal associated
244 with PCOS. We validated the association of previously reported genes *ERBB4* and *DENND1A1*,
245 with rs113168128 being an independent signal in *ERBB4*. We also identified rs144248326, an
246 intronic variant of *WWTR1*, as a novel signal close to genome-wide significance level. In further
247 association tests of the top 3 SNPs with each criterion in the PCOS algorithm, we found that
248 rs17186366 was associated with polycystic and hyperandrogenism, while rs11316812 and
249 rs144248326 were strongly associated with oligomenorrhea or infertility. The top SNPs are

250 associated with traits related to the nervous system and mental/behavior disorders in the
251 PheWAS queries.

252 Results and Implications

253 During case identification, we observed similar proportions in each criterion used in the
254 algorithm for the eMERGE and BioVU data, but different from the MyCode data. Around 40% of
255 the patients had polycystic ovaries in the eMERGE and BioVU cohorts, versus 88% in the
256 Geisinger cohort. This may be due to the integration of information from the patients' problem
257 list at Geisinger or a difference in clinical practice between the systems.

258 We identified a novel genome-wide significant variant rs17186366 near *SOD2*, which was
259 associated with the polycystic ovaries and hyperandrogenism. *SOD2* encodes Superoxide
260 Dismutase 2, a mitochondrial protein which converts superoxide byproducts of oxidative
261 phosphorylation to hydrogen peroxide and diatomic oxygen. Recently, A16V (rs4880) in *SOD2*
262 was found to be associated with PCOS, serum luteinizing hormone (LH) levels, and the ratio of
263 LH to follicle-stimulating hormone in Han-Chinese women³⁵. rs17186366 was also found to be
264 associated with menarche in the UKB GWAS results. One retrospective study showed early or
265 late menarche were more likely to be seen in women with PCOS³⁶.

266 The association of *ERBB4* with PCOS has been validated in both European and Han-Chinese^{13, 16,}
267 ³⁷. In this study, we identified an intronic variant rs11316812 of *ERBB4* that is not in LD with the
268 known variants at close to genome-wide significant level ($P=5.2 \times 10^{-8}$). *ERBB4*, also known as
269 human epidermal growth factor receptor 4 (HER4), belongs to the EGFR family which includes
270 *ERBB1*, *ERBB2/HER2* and *ERBB3/HER3*. Other than PCOS, variants of *ERBB4* have been

271 associated with ovarian cancer³⁸ and schizophrenia³⁹. ERBB4 can be stimulated by its ligands
272 and activate the EGFR signaling, which is critical for LH-induced steroidogenesis which
273 promotes follicular maturation^{40,41}. ERBB4 signaling is also involved in luteal growth⁴². ERBB4
274 is highly expressed in the nerves system according to GTEx datasets, including in the
275 hypothalamus and pituitary, the two important organs in the hypothalamus-pituitary-ovary-
276 adrenal (HPOA) axis. These findings suggest a disturbance in the control mechanisms of the
277 HPOA axis in PCOS.

278 *WWTR1*, the third gene with strong association, encodes for TAZ (transcriptional co-activator
279 with PDZ-binding motif). TAZ also contains a WW domain as the Yes-associated protein (YAP1)
280⁴³ and is another gene that was associated with PCOS^{12, 13, 16}. *WWTR1* and YAP1 are two key
281 molecules of the Hippo signaling pathway, and their expression were significantly altered in
282 PCOS tissues⁴⁴. Insulin resistance affects 50-70% of women with PCOS⁴⁵. *WWTR1* and YAP1 can
283 also regulate insulin resistance. The inhibition of *WWTR1*/YAP1 in combination with metformin
284 can completely inhibit the effect of insulin⁴⁶. Our findings provide a possible genetic link
285 between PCOS and the Hippo pathway, suggesting potential pharmaceutical Hippo-targeted
286 interventions for treatment of PCOS with insulin resistance. Interestingly, ERBB4 can also
287 interact with WW domains in YAP1 through the PPxY motif⁴⁷. The ERBB4-YAP1-*WWTR1*
288 interaction network indicates the Hippo and EGFR signaling contribute to the multifactorial
289 etiology of PCOS.

290 Epidemiologic studies found that women with PCOS have a higher risk for depression, bipolar
291 disorder, anxiety, and eating disorders⁴⁵. In our study, we found moderate associations for the
292 top three variants with mental/behavioral disorders including depression. When queried, the

293 PheWAS and GWAS catalog through Open Targets Genetics, suggests shared genetic
294 predisposition between PCOS and mental disorders. Also, the differences in association
295 patterns across the PCOS criteria likely reflect the complex biology of PCOS and support the
296 epidemiological observations of heterogeneity in the phenotype.

297 [Strengths and Limitations](#)

298 The major strength of this study is the same phenotyping algorithm was applied through
299 different systems to ensure the homogeneity. Our PCOS algorithm was based on the Rotterdam
300 criteria and thus was broader than the NIH or AE-PCOS criteria, enabling us to identify more
301 cases than the ICD code-based method. This would improve the number of cases and avoid the
302 cases that would be identified as controls otherwise, thus increases the specificity of the study.
303 However, the algorithm assumes the same evaluation and coding practices at each healthcare
304 system and its sensitivity and specificity for case identification was not validated by chart
305 review. A potential limitation is case definition of testosterone laboratory measurement
306 (criteria2c) was unavailable in the eMERGE data. However, based on MyCode, all the patients
307 identified by this criterion also met the other two criteria. Thus, absence of this information
308 may not affect the final total case number. Also, the EHR data on luteinizing hormone and
309 follicle-stimulating hormone serum levels were limited thus we were not able to test the
310 associations with the top variants. The sample size is not large enough to identify variants with
311 small effect or low minor allele frequency, especially because PCOS is a heterogenous disorder
312 with complex etiology and combinations of clinical symptoms. For our study, we have about
313 80% power for the top 3 variants. The study cohorts are mainly European-decent individuals,
314 with a very small proportion of African Americans and other race/ethnic groups. Although the

315 variant in *ERBB4* showed the same direction of effects and nominal significance ($P < 0.05$) in the
316 African American, it has a higher minor allele frequency than in the European population.
317 Future studies focusing on minorities are necessary.

318 Conclusions

319 Through a discovery-validation GWAS on PCOS cases and controls identified from EHR using an
320 algorithm based on Rotterdam criteria, we validated the association with *ERBB4*. We also
321 identified novel association with *SOD2* and *WWTR1*. Our findings highlighted the role of EGFR
322 and Hippo signaling in the disturbance of metabolic and HPOA axis in PCOS etiology.

323 Acknowledgements

324 The authors thank Christina M. Yule and Sara J. Kwiecien at Geisinger, and Brittany City at the
325 eMERGE network coordinating center. The authors express their gratitude to Drs. Cecilia
326 Lindgren, John Perry, and Corrine Kolka Welt at the International PCOS Consortium for their
327 support. The authors would like to thank Ilene Ladd for English editing.

328 Authors contribution

329 YZ, KH, AJ, and ML designed the study; KH and DH developed the algorithm and applied to
330 Geisinger EHR; YZ performed the phenotyping for eMERGE data and discovery GWAS; Jacob
331 performed the phenotyping and replication using BioVU data; FD extracted the PCOS
332 consortium data; YZ drafted the manuscript; KH, AJ, ML, SP and MSW did critical review; NJ, LD,
333 AR, GJ, IS, DVE, PS, EA and BH contributed the eMERGE data and provided critical review; All
334 authors approved the final manuscript.

335 Conflict of interest

336 The authors report no conflict of interest.

337 Funding disclosure

338 MyCode® was funded by Geisinger and Regeneron Genomics Center; the eMERGE III was
339 funded by NIH U01HG8679 (Geisinger Clinic). The funding sources was not involved in the
340 interpretation of the result or which journal to submit.

341

342 References

- 343 1. HART R, HICKEY M, FRANKS S. Definitions, prevalence and symptoms of polycystic
344 ovaries and polycystic ovary syndrome. *Best Pract Res Clin Obstet Gynaecol*
345 2004;18:671-83.
- 346 2. ZAWADZKI JK DA. *Diagnostic criteria for polycystic ovary syndrome: towards a*
347 *rational approach*. Boston: Blackwell Scientific Publications; Number of pages.
- 348 3. ROTTERDAM EA-SPCWG. Revised 2003 consensus on diagnostic criteria and long-
349 term health risks related to polycystic ovary syndrome. *Fertil Steril* 2004;81:19-25.
- 350 4. AZZIZ R, CARMINA E, DEWAILLY D, et al. Positions statement: criteria for defining
351 polycystic ovary syndrome as a predominantly hyperandrogenic syndrome: an
352 Androgen Excess Society guideline. *J Clin Endocrinol Metab* 2006;91:4237-45.
- 353 5. MARCH WA, MOORE VM, WILLSON KJ, PHILLIPS DI, NORMAN RJ, DAVIES MJ. The
354 prevalence of polycystic ovary syndrome in a community sample assessed under
355 contrasting diagnostic criteria. *Hum Reprod* 2010;25:544-51.
- 356 6. YILDIZ BO, BOZDAG G, YAPICI Z, ESINLER I, YARALI H. Prevalence, phenotype and
357 cardiometabolic risk of polycystic ovary syndrome under different diagnostic
358 criteria. *Hum Reprod* 2012;27:3067-73.
- 359 7. VINK JM, SADRZADEH S, LAMBALK CB, BOOMSMA DI. Heritability of polycystic ovary
360 syndrome in a Dutch twin-family study. *J Clin Endocrinol Metab* 2006;91:2100-4.
- 361 8. JAHANFAR S, EDEN JA, NGUYEN T, WANG XL, WILCKEN DE. A twin study of polycystic
362 ovary syndrome and lipids. *Gynecol Endocrinol* 1997;11:111-7.
- 363 9. JAHANFAR S, EDEN JA, WARREN P, SEPPALA M, NGUYEN TV. A twin study of polycystic
364 ovary syndrome. *Fertil Steril* 1995;63:478-86.
- 365 10. CHEN ZJ, ZHAO H, HE L, et al. Genome-wide association study identifies
366 susceptibility loci for polycystic ovary syndrome on chromosome 2p16.3, 2p21
367 and 9q33.3. *Nat Genet* 2011;43:55-9.

- 368 11. HWANG JY, LEE EJ, JIN GO M, et al. Genome-wide association study identifies GYS2
369 as a novel genetic factor for polycystic ovary syndrome through obesity-related
370 condition. *J Hum Genet* 2012;57:660-4.
- 371 12. SHI Y, ZHAO H, SHI Y, et al. Genome-wide association study identifies eight new risk
372 loci for polycystic ovary syndrome. *Nat Genet* 2012;44:1020-5.
- 373 13. DAY FR, HINDS DA, TUNG JY, et al. Causal mechanisms and balancing selection
374 inferred from genetic associations with polycystic ovary syndrome. *Nat Commun*
375 2015;6:8464.
- 376 14. HAYES MG, URBANEK M, EHRMANN DA, et al. Genome-wide association of polycystic
377 ovary syndrome implicates alterations in gonadotropin secretion in European
378 ancestry populations. *Nat Commun* 2015;6:7502.
- 379 15. LEE H, OH JY, SUNG YA, et al. Genome-wide association study identified new
380 susceptibility loci for polycystic ovary syndrome. *Hum Reprod* 2015;30:723-31.
- 381 16. DAY F, KARADERI T, JONES MR, et al. Large-scale genome-wide meta-analysis of
382 polycystic ovary syndrome suggests shared genetic architecture for different
383 diagnosis criteria. *PLoS Genet* 2018;14:e1007813.
- 384 17. ABUL-HUSN NS, KENNY EE. Personalized Medicine and the Power of Electronic
385 Health Records. *Cell* 2019;177:58-69.
- 386 18. CAREY DJ, FETTEROLF SN, DAVIS FD, et al. The Geisinger MyCode community health
387 initiative: an electronic health record-linked biobank for precision medicine
388 research. *Genet Med* 2016;18:906-13.
- 389 19. RODEN DM, PULLEY JM, BASFORD MA, et al. Development of a large-scale de-
390 identified DNA biobank to enable personalized medicine. *Clin Pharmacol Ther*
391 2008;84:362-9.
- 392 20. PULLEY J, CLAYTON E, BERNARD GR, RODEN DM, MASYS DR. Principles of human
393 subjects protections applied in an opt-out, de-identified biobank. *Clin Transl Sci*
394 2010;3:42-8.
- 395 21. McCARTY CA, CHISHOLM RL, CHUTE CG, et al. The eMERGE Network: a consortium of
396 biorepositories linked to electronic medical records data for conducting genomic
397 studies. *BMC Med Genomics* 2011;4:13.
- 398 22. KHO AN, PACHECO JA, PEISSIG PL, et al. Electronic medical records for genetic
399 research: results of the eMERGE consortium. *Sci Transl Med* 2011;3:79re1.
- 400 23. GOTTESMAN O, KUIVANIEMI H, TROMP G, et al. The Electronic Medical Records and
401 Genomics (eMERGE) Network: past, present, and future. *Genet Med* 2013;15:761-
402 71.
- 403 24. STANAWAY IB, HALL TO, ROSENTHAL EA, et al. The eMERGE genotype set of 83,717
404 subjects imputed to ~40 million variants genome wide and association with the
405 herpes zoster medical record phenotype. *Genet Epidemiol* 2019;43:63-81.

- 406 25. NEWTON KM, PEISSIG PL, KHO AN, et al. Validation of electronic medical record-
407 based phenotyping algorithms: results and lessons learned from the eMERGE
408 network. *J Am Med Inform Assoc* 2013;20:e147-54.
- 409 26. DEWEY FE, MURRAY MF, OVERTON JD, et al. Distribution and clinical impact of
410 functional variants in 50,726 whole-exome sequences from the DiscovEHR study.
411 *Science* 2016;354.
- 412 27. ZHANG Y, POLER SM, LI J, et al. Dissecting genetic factors affecting phenylephrine
413 infusion rates during anesthesia: a genome-wide association study employing
414 EHR data. *BMC Med* 2019;17:168.
- 415 28. WINKLER TW, DAY FR, CROTEAU-CHONKA DC, et al. Quality control and conduct of
416 genome-wide association meta-analyses. *Nat Protoc* 2014;9:1192-212.
- 417 29. WILLER CJ, LI Y, ABECASIS GR. METAL: fast and efficient meta-analysis of
418 genomewide association scans. *Bioinformatics* 2010;26:2190-1.
- 419 30. CHANG CC, CHOW CC, TELLIER LC, VATTIKUTI S, PURCELL SM, LEE JJ. Second-generation
420 PLINK: rising to the challenge of larger and richer datasets. *Gigascience* 2015;4:7.
- 421 31. McLAREN W, GIL L, HUNT SE, et al. The Ensembl Variant Effect Predictor. *Genome*
422 *Biol* 2016;17:122.
- 423 32. WATANABE K, TASKESEN E, VAN BOCHOVEN A, POSTHUMA D. Functional mapping and
424 annotation of genetic associations with FUMA. *Nat Commun* 2017;8:1826.
- 425 33. CARVALHO-SILVA D, PIERLEONI A, PIGNATELLI M, et al. Open Targets Platform: new
426 developments and updates two years on. *Nucleic Acids Res* 2019;47:D1056-D65.
- 427 34. NAGEL M, JANSEN PR, STRINGER S, et al. Meta-analysis of genome-wide association
428 studies for neuroticism in 449,484 individuals identifies novel genetic loci and
429 pathways. *Nat Genet* 2018;50:920-27.
- 430 35. LIU Q, LIU H, BAI H, et al. Association of SOD2 A16V and PON2 S311C
431 polymorphisms with polycystic ovary syndrome in Chinese women. *J Endocrinol*
432 *Invest* 2019.
- 433 36. CARROLL J, SAXENA R, WELT CK. Environmental and genetic factors influence age at
434 menarche in women with polycystic ovary syndrome. *J Pediatr Endocrinol Metab*
435 2012;25:459-66.
- 436 37. PENG Y, ZHANG W, YANG P, et al. ERBB4 Confers Risk for Polycystic Ovary Syndrome
437 in Han Chinese. *Sci Rep* 2017;7:42000.
- 438 38. WEI P, LI L, ZHANG Z, ZHANG W, LIU M, SHENG X. A genetic variant of miR-335 binding
439 site in the ERBB4 3'-UTR is associated with prognosis of ovary cancer. *J Cell*
440 *Biochem* 2018;119:5135-42.
- 441 39. SILBERBERG G, DARVASI A, PINKAS-KRAMARSKI R, NAVON R. The involvement of ErbB4
442 with schizophrenia: association and expression studies. *Am J Med Genet B*
443 *Neuropsychiatr Genet* 2006;141B:142-8.
- 444 40. PARK JY, SU YQ, ARIGA M, LAW E, JIN SL, CONTI M. EGF-like growth factors as
445 mediators of LH action in the ovulatory follicle. *Science* 2004;303:682-4.

- 446 41. JAMNONGJIT M, GILL A, HAMMES SR. Epidermal growth factor receptor signaling is
447 required for normal ovarian steroidogenesis and oocyte maturation. *Proc Natl*
448 *Acad Sci U S A* 2005;102:16257-62.
- 449 42. AKAYAMA Y, TAKEKIDA S, OHARA N, et al. Gene expression and immunolocalization of
450 heparin-binding epidermal growth factor-like growth factor and human
451 epidermal growth factor receptors in human corpus luteum. *Hum Reprod*
452 2005;20:2708-14.
- 453 43. KANAI F, MARIGNANI PA, SARBASSOVA D, et al. TAZ: a novel transcriptional co-
454 activator regulated by interactions with 14-3-3 and PDZ domain proteins. *EMBO J*
455 2000;19:6778-91.
- 456 44. MAAS K, MIRABAL S, PENZIAS A, SWEETNAM PM, EGGAN KC, SAKKAS D. Hippo signaling in
457 the ovary and polycystic ovarian syndrome. *J Assist Reprod Genet* 2018;35:1763-
458 71.
- 459 45. SIRMANS SM, PATE KA. Epidemiology, diagnosis, and management of polycystic
460 ovary syndrome. *Clin Epidemiol* 2013;6:1-13.
- 461 46. WANG C, JEONG K, JIANG H, et al. YAP/TAZ regulates the insulin signaling via IRS1/2
462 in endometrial cancer. *Am J Cancer Res* 2016;6:996-1010.
- 463 47. HASKINS JW, NGUYEN DX, STERN DF. Neuregulin 1-activated ERBB4 interacts with
464 YAP to induce Hippo pathway target genes and promote cell migration. *Sci Signal*
465 2014;7:ra116.
- 466
- 467

468 **Tables**

469 **Table 1. Characteristics of PCOS cases and controls from discovery and replication cohorts**

Cohorts	Discovery						Replication	
	MyCode Phase I		MyCode Phase II		eMERGE Phase III		BioVU	
	Case	Control	Case	Control	Case	Control	Case	Control
Number	1,141	18,788	594	9,024	1,007	23,626	253	2,161
- Polycystic (C1)	1,011 (88.6%)	/	528 (88.9%)	/	390 (38.7%)	/	107(42.1%)	/
- Hyperandrogenic (C2)	773 (67.8%)	/	385 (64.8%)	/	841 (83.5%)	/	209 (82.3%)	/
- Reproductive (C3)	1,120 (98.2%)	/	583 (98.1%)	/	945 (93.8%)	/	240 (94.5%)	/
Age, Mean (SD)	39.5 (8.4)	64.8(13.5)	37.7 (8.7)	60.9 (12.4)	33.1 (7.0)	70.7 (16.1)	35.4 (7.2)	41.4(2.3)
BMI, Mean (SD)	35.0 (9.8)	31.1 (7.9)	34.3 (9.5)	31.0 (7.6)	30.4 (8.4)	30.0 (7.1)	29.5 (8.8)	28.2(8.1)

470

471 **Table 2: Association of the top 3 SNPs with each PCOS criterion**

SNP	Criterion	MyCode Phase I		MyCode Phase II		eMERGE III		Meta	
		OR	P	OR	P	OR	P	OR [CI95]	P
rs17186366 <i>SOD2</i>	Polycystic	1.44	0.0004	1.54	0.0026	1.38	0.0295	1.45[1.25,1.67]	3.36x10 ⁻⁷
	Hyperandrogenic	1.49	0.0003	1.48	0.0163	1.27	0.0258	1.39[1.21,1.59]	1.81x10 ⁻⁶
	Reproductive	1.31	0.0052	1.26	0.0992	1.27	0.0213	1.28[1.13,1.45]	7.23x10 ⁻⁵
rs113168128 <i>ERBB4</i>	Polycystic	2.08	0.0004	2.19	0.0148	1.34	0.2024	1.79[1.36,2.35]	2.85x10 ⁻⁵
	Hyperandrogenic	1.39	0.1818	2.27	0.0186	1.59	0.0033	1.60[1.26,2.04]	1.31x10 ⁻⁴
	Reproductive	1.88	0.0013	2.32	0.0040	1.62	0.0012	1.79[1.44,2.22]	1.40x10 ⁻⁷
rs144248326 <i>WWTR1</i>	Polycystic	1.91	0.0276	2.71	0.0240	0.91	0.8343	1.74[1.15,2.65]	9.43x10 ⁻³
	Hyperandrogenic	1.90	0.0406	2.50	0.0625	2.23	0.0002	2.16[1.55,3.00]	4.94x10 ⁻⁶
	Reproductive	2.18	0.0032	3.49	0.0011	1.87	0.0052	2.19[1.61,2.97]	5.31x10 ⁻⁷

472

473

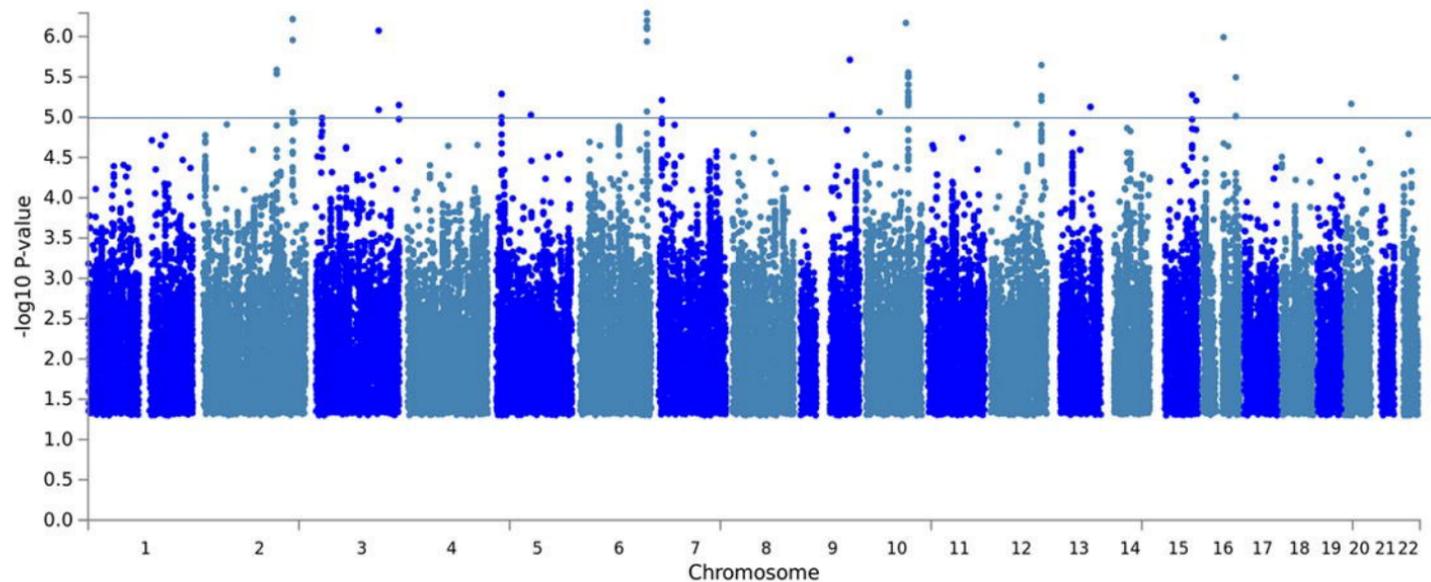
474 **Figure legends**

475 **Figure 1: Flowchart for study design.** The PCOS algorithm was first developed to identify cases
 476 and controls from the EHR based on Rotterdam criteria, and then applied to Geisinger patients,
 477 and the eMERGE cohort. Case-control GWAS were then conducted for the three cohorts with
 478 genetics data followed by fixed effect inverse-variance meta-analyses. Variants with P<1e-6
 479 were then validated in BioVU samples using the same phenotype algorithm and genetics
 480 analyses protocol for in meta-analysis and were queried in summary statistics from the PCOS
 481 consortium. Association with each criterion in the PCOS algorithm were further tested for these
 482 variants.

483 **Figure 2: GWAS analyses of PCOS and functional evaluation.** (A) Manhattan plot for the meta-
484 analysis of the Geisinger MyCode Phase I, II and the eMERGE Phase III cohorts. (B) The top 3
485 associated variants with P value < 1e-6 in discovery meta-analysis and their replication and final
486 meta-analysis. (C) The protein-protein interaction (PPI) network for ERBB4, WWTR1 and YAP1
487 using STRING. Only high confidence interactions were shown (confidence score >=0.7).

488

A



B

SNP	CHR:BP	EA/RA	EAF	Nearest Gene	OR [CI95]	P-discov	P-rep	P-meta
rs17186366	6:159898261	C/T	0.0836	<i>LOC101929142</i> <i>/SOD2</i>	1.37 [1.23,1.54]	7.78×10^{-7}	0.0088	2.80×10^{-8}
rs113168128 *	2:212291772	A/G	0.0206	<i>ERBB4</i>	1.72 [1.42, 2.10]	6.07×10^{-7}	0.064	5.20×10^{-8}
rs144248326	3:149319873	C/T	0.0103	<i>WWTR1</i>	2.13 [1.56, 2.86]	8.45×10^{-7}	/	8.45×10^{-7}

* $R^2 = 0.001$ with reported SNP rs2178575 in CEU 1000 Genome data. Abbreviations: EA: effect allele; RA: reference allele; OR: odds ratio; CI: confidence interval; P-discov: P value of discovery meta-analysis; P-rep: P value of replication using BioVU samples; P-meta: P value of final meta-analyses .

C

