

1 **Title: Predicting circulating CA125 levels among healthy premenopausal women**

2

3 **Authors and affiliations:** Naoko Sasamoto<sup>1,2</sup>, Ana Babic,<sup>2,3</sup> Bernard A. Rosner<sup>2,4</sup>, Renée T.  
4 Fortner<sup>5</sup>, Allison F. Vitonis<sup>1</sup>, Hidemi Yamamoto<sup>6</sup>, Raina N. Fichorova<sup>2,6</sup>, Anne Tjønneland<sup>7</sup>,  
5 Louise Hansen<sup>7</sup>, Kim Overvad<sup>8</sup>, Marina Kvaskoff<sup>9,10</sup>, Agnès Fournier<sup>9,10</sup>, Francesca Romana  
6 Mancini<sup>9,10</sup>, Heiner Boeing<sup>11</sup>, Antonia Trichopoulou<sup>12,13</sup>, Eleni Peppas<sup>12</sup>, Anna Karakatsani<sup>12,14</sup>,  
7 Domenico Palli<sup>15</sup>, Valeria Pala<sup>16</sup>, Amalia Mattiello<sup>17</sup>, Rosario Tumino<sup>18</sup>, Chiara Grasso<sup>19</sup>, N.  
8 Charlotte Onland-Moret<sup>20</sup>, Elisabete Weiderpass<sup>21-24</sup>, J. Ramón Quirós<sup>25</sup>, Leila Lujan-Barroso<sup>26</sup>,  
9 Miguel Rodríguez-Barranco<sup>27,28</sup>, Sandra Colorado-Yohar<sup>29,30</sup>, Aurelio Barricarte<sup>31,32</sup>, Miren  
10 Dorronsoro<sup>33</sup>, Annika Idahl<sup>34,35</sup>, Eva Lundin<sup>36</sup>, Hanna Sartor<sup>37,38</sup>, Kay-Tee Khaw<sup>39</sup>, Timothy J  
11 Key<sup>40</sup>, David Muller<sup>41</sup>, Elio Riboli<sup>42</sup>, Marc Gunter<sup>42</sup>, Laure Dossus<sup>42</sup>, Rudolf Kaaks<sup>5</sup>, Daniel W.  
12 Cramer<sup>1,2</sup>, Shelley S. Tworoger<sup>43</sup>, Kathryn L. Terry<sup>1,2</sup>

- 13 1. Obstetrics and Gynecology Epidemiology Center, Brigham and Women's Hospital,  
14 Boston, Massachusetts
- 15 2. Harvard Medical School, Boston, Massachusetts
- 16 3. Department of Medical Oncology, Dana-Farber Cancer Institute, Boston, Massachusetts
- 17 4. Channing Division of Network Medicine, Department of Medicine, Brigham and Women's  
18 Hospital, Boston, Massachusetts
- 19 5. Division of Cancer Epidemiology, German Cancer Research Center (DKFZ) Heidelberg,  
20 Germany.
- 21 6. Laboratory of Genital Tract Biology, Department of Obstetrics, Gynecology and  
22 Reproductive Biology, Brigham and Women's Hospital, Boston, Massachusetts
- 23 7. Diet, Genes and Environment, Danish Cancer Society Research Center, Copenhagen,  
24 Denmark
- 25 8. Department of Public Health, Section for Epidemiology, Aarhus University, Aarhus,  
26 Denmark
- 27 9. CESP, Fac. de médecine - Univ. Paris-Sud, Fac. de médecine - UVSQ, INSERM,  
28 Université Paris-Saclay, Villejuif, France
- 29 10. Gustave Roussy, Villejuif, France
- 30 11. Department of Epidemiology, German Institute of Human Nutrition Potsdam-Rehbruecke,  
31 Nuthetal, Germany

- 32 12. Hellenic Health Foundation, Athens, Greece
- 33 13. WHO Collaborating Center for Nutrition and Health, Unit of Nutritional Epidemiology  
34 and Nutrition in Public Health, Dept. of Hygiene, Epidemiology and Medical Statistics,  
35 School of Medicine, National and Kapodistrian University of Athens, Greece
- 36 14. 2nd Pulmonary Medicine Department, School of Medicine,  
37 National and Kapodistrian University of Athens, "ATTIKON" University Hospital,  
38 Haidari, Greece
- 39 15. Head, Cancer Risk Factors and Life-Style Epidemiology Unit  
40 Institute for Cancer Research, Prevention and Clinical Network - ISPRO,  
41 Florence, Italy
- 42 16. Epidemiology and Prevention Unit, Fondazione IRCCS Istituto Nazionale dei Tumori di  
43 Milano, Milano, Italy
- 44 17. Dipartimento Di Medicina Clinica E Chirurgia, Federico II University, Naples, Italy
- 45 18. Cancer Registry and Histopathology Department, "Civic - M.P. Arezzo" Hospital, ASP  
46 Ragusa, Italy
- 47 19. Unit of Cancer Epidemiology– CeRMS, Department of Medical Sciences, University of  
48 Turin, Turin, Italy
- 49 20. Department of Epidemiology, Julius Center for Health Sciences and Primary Care,  
50 University Medical Center Utrecht
- 51 21. Department of Research, Cancer Registry of Norway, Institute of Population-Based  
52 Cancer Research, Oslo, Norway
- 53 22. Department of Medical Epidemiology and Biostatistics, Karolinska Institutet, Stockholm,  
54 Sweden
- 55 23. Genetic Epidemiology Group, Folkhälsan Research Center and Faculty of Medicine,  
56 Helsinki University, Helsinki, Finland
- 57 24. Department of Community Medicine, University of Tromsø , The Arctic University of  
58 Norway, Tromsø, Norway
- 59 25. Public Health Directorate, Asturias, Spain
- 60 26. Unit of Nutrition and Cancer. Cancer Epidemiology Research Program, Catalan Institute  
61 of Oncology (ICO-IDIBELL), L'Hospitalet de Llobregat, Barcelona, Spain

- 62 27. Escuela Andaluza de Salud Pública. Instituto de Investigación Biosanitaria  
63 ibs.GRANADA. Hospitales Universitarios de Granada/Universidad de Granada,  
64 Granada, Spain
- 65 28. CIBER de Epidemiología y Salud Pública (CIBERESP), Madrid, Spain
- 66 29. Department of Epidemiology, Murcia Health Council, IMIB-Arrixaca, Spain.
- 67 30. Research Group on Demography and Health, National Faculty of Public Health,  
68 University of Antioquia, Medellín, Colombia
- 69 31. Navarra Public Health Institute, Pamplona, Spain Navarra Institute for Health Research  
70 (IdiSNA) Pamplona, Spain
- 71 32. CIBER Epidemiology and Public Health CIBERESP, Madrid. Spain
- 72 33. Public Health Direction and Biodonostia Research Institute and Ciberesp, Basque  
73 Regional Health Department, San Sebastian. Spain
- 74 34. Department of Clinical Sciences, Obstetrics and Gynecology, Umeå University, Umeå,  
75 Sweden
- 76 35. Department of Public Health and Clinical Medicine, Nutritional Research, Umeå  
77 University, Umeå, Sweden
- 78 36. Department of Medical Biosciences, Pathology, Umeå University, Umeå, Sweden
- 79 37. Department of Medical Imaging and Physiology, Skåne University Hospital, Lund,  
80 Sweden
- 81 38. Department of Translational Medicine, Lund University, Sweden
- 82 39. Cancer Epidemiology Unit, University of Cambridge, Cambridge, United Kingdom
- 83 40. Cancer Epidemiology Unit, Nuffield Department of Population Health, University of  
84 Oxford
- 85 41. Department of Epidemiology and Biostatistics, School of Public Health, Imperial College  
86 London, United Kingdom
- 87 42. International Agency for Research on Cancer, Lyon, France
- 88 43. Department of Cancer Epidemiology, Moffitt Cancer Center, Tampa, Florida

89  
90 **Running title:** CA125 prediction model in premenopausal women

91 **Keywords:** CA125, prediction model, premenopausal, ovarian cancer, screening

92

93 **Additional Information:**

94 **Financial support:** This study was funded by National Institutes of Health R01 CA193965, R01  
95 CA 158119, and R35 CA197605. The coordination of EPIC is financially supported by the  
96 European Commission (DG-SANCO) and the International Agency for Research on Cancer. The  
97 national cohorts are supported by Danish Cancer Society (Denmark); Ligue Contre le Cancer,  
98 Institut Gustave Roussy, Mutuelle Générale de l'Éducation Nationale, Institut National de la Santé  
99 et de la Recherche Médicale (INSERM) (France); German Cancer Aid, German Cancer Research  
100 Center (DKFZ), Federal Ministry of Education and Research (BMBF), Deutsche Krebshilfe,  
101 Deutsches Krebsforschungszentrum and Federal Ministry of Education and Research (Germany);  
102 the Hellenic Health Foundation (Greece); Associazione Italiana per la Ricerca sul Cancro-AIRC-  
103 Italy and National Research Council (Italy); Dutch Ministry of Public Health, Welfare and Sports  
104 (VWS), Netherlands Cancer Registry (NKR), LK Research Funds, Dutch Prevention Funds, Dutch  
105 ZON (Zorg Onderzoek Nederland), World Cancer Research Fund (WCRF), Statistics Netherlands  
106 (The Netherlands); ERC-2009-AdG 232997 and Nordforsk, Nordic Centre of Excellence  
107 programme on Food, Nutrition and Health (Norway); Health Research Fund (FIS), PI13/00061 to  
108 Granada; , PI13/01162 to EPIC-Murcia), Regional Governments of Andalucía, Asturias, Basque  
109 Country, Murcia and Navarra, ISCIII RETIC (RD06/0020) (Spain); Swedish Cancer Society,  
110 Swedish Research Council and County Councils of Skåne and Västerbotten, The Cancer Research  
111 Foundation of Northern Sweden (Sweden); Cancer Research UK (14136 to EPIC-Norfolk;  
112 C570/A16491 and C8221/A19170 to EPIC-Oxford), Medical Research Council (1000143 to  
113 EPIC-Norfolk, MR/M012190/1 to EPIC-Oxford) (United Kingdom).

114

115

116

117 **Corresponding author:**

118 Naoko Sasamoto

119 Obstetrics and Gynecology Epidemiology Center

120 Brigham and Women's Hospital

121 221 Longwood Avenue, Boston, MA 02115

122 Phone: 617-732-4895, Fax: 617-732-4899

123 Email: [nsasmoto@bwh.harvard.edu](mailto:nsasmoto@bwh.harvard.edu)

124

125 **Conflict of interest:** The authors declare no potential conflicts of interest.

126 **Abstract word count:** 250 words

127 **Manuscript word count:** 3,899 words

128 **Number of figures and tables:** 3 Tables, 3 Figures, 4 Supplementary Tables

129

130

131

132

133

134

135

136

137

138

139

140

141

142

143

144

145 **Abstract**

146 Background: CA125 is the most promising ovarian cancer screening biomarker to date. Multiple  
147 studies reported CA125 levels vary by personal characteristics, which could inform personalized  
148 CA125 thresholds. However, this has not been well described in premenopausal women.

149 Methods: We evaluated predictors of CA125 levels among 815 premenopausal women from the  
150 New England Case Control Study (NEC). We developed linear and dichotomous ( $\geq 35$  U/ mL)  
151 CA125 prediction models and externally validated an abridged model restricting to available  
152 predictors among 473 premenopausal women in the European Prospective Investigation into  
153 Cancer and Nutrition Study (EPIC).

154 Results: The final linear CA125 prediction model included age, race, tubal ligation, endometriosis,  
155 menstrual phase at blood draw, and fibroids, which explained 7% of the total variance of CA125.  
156 The correlation between observed and predicted CA125 levels based on the abridged model  
157 (including age, race, and menstrual phase at blood draw) had similar correlation coefficients in  
158 NEC( $r=0.22$ ) and in EPIC( $r=0.22$ ). The dichotomous CA125 prediction model included age, tubal  
159 ligation, endometriosis, prior personal cancer diagnosis, family history of ovarian cancer, number  
160 of miscarriages, menstrual phase at blood draw and smoking status with AUC of 0.83. The  
161 abridged dichotomous model (including age, number of miscarriages, menstrual phase at blood  
162 draw, and smoking status) showed similar AUCs in NEC(0.73) and in EPIC(0.78).

163 Conclusions: We identified a combination of factors associated with CA125 levels in  
164 premenopausal women.

165 Impact: Our model could be valuable in identifying healthy women likely to have elevated CA125  
166 and consequently improve its specificity for ovarian cancer screening.

167

168

169

170

171 **Introduction**

172 Ovarian cancer is the eighth leading cause of cancer death in 2012 with 151,900 deaths worldwide  
173 due, in part, to lack of specific symptoms leading to diagnosis at late stage when prognosis is  
174 poor(1,2). More than 80% of ovarian cancer patients have elevated cancer antigen 125 (CA125),  
175 a membrane bound glycosylated mucin (MUC16), which is used clinically as a prognostic  
176 biomarker and to monitor response to therapy(3,4). However, results from two large randomized  
177 screening trials in primarily postmenopausal women using transvaginal ultrasound and CA125  
178 (either using 35 U/ml as a cutoff, or the risk of ovarian cancer algorithm (ROCA)) showed no  
179 clinically significant benefit(5,6). MUC16 is expressed on a variety of tissues, including the lung,  
180 pancreas, stomach, liver, endometrium, and breast, and levels vary between individuals based on  
181 demographic, reproductive and lifestyle characteristics(7-11). Therefore, identifying personal  
182 characteristics that are associated with CA125 levels could be used to create personalized  
183 thresholds for CA125 instead of a single 35 U/mL cutoff, thereby improving the interpretation of  
184 measured CA125 and its performance as a screening biomarker and ultimately leading to decreased  
185 ovarian cancer mortality.

186 However, prior studies examining factors associated with CA125 have focused on postmenopausal  
187 women(7,8). Thus, we evaluated factors associated with CA125 in premenopausal women and  
188 developed and validated CA125 prediction models (linear and dichotomous) among  
189 premenopausal women without ovarian cancer from the New England Case-Control Study and the  
190 European Prospective Investigation into Nutrition and Cancer study.

191

192 **Materials and Methods**

193 **Study population**

194 The New England Case Control Study (NEC) is a population-based case-control study of ovarian  
195 cancer with 2,100 population-based controls enrolled from New Hampshire and eastern  
196 Massachusetts over the three study phases (1992-97, 1998-2002, 2003-2008). Details on the study  
197 design have been described previously(12-14). Briefly, controls were identified using random-  
198 digit dialing, town book selection, and drivers' license lists and frequency matched on age and

199 state of residence. Approximately half (54%) of the eligible controls that were contacted agreed to  
200 participate. We restricted the study population to controls (n=2,100) and excluded women without  
201 CA125 measurements (n=96), women postmenopausal at time of blood draw (n=1,130), women  
202 who had hysterectomy due to unknown menopausal status (n=30), women who were pregnant or  
203 breastfeeding at time of blood draw (n=25), and women with extreme CA125 values ranging from  
204 115.3U/ mL to 411.7 U/mL (n=4) identified based on the generalized extreme studentized deviated  
205 many-outlier detection approach applied to log-transformed values(15). In sum, our analysis  
206 included 815 premenopausal NEC controls.

207 The European Prospective Investigation into Cancer and Nutrition (EPIC) study is a multicenter  
208 prospective cohort including participants from ten Western European countries developed to  
209 evaluate the association between nutrition and cancer. Briefly, 519,978 participants (366,521  
210 women) were enrolled between 1991 and 1998 across 23 research centers. Details on the study  
211 design have been described previously(16). A nested case-control study of ovarian cancer was  
212 designed within the cohort(17). For each ovarian cancer case, up to four controls were randomly  
213 selected using incidence density sampling for a total of 1,939 controls(17). We excluded women  
214 without CA125 measurements (n=12), women who were either postmenopausal (n=1,416), or had  
215 a hysterectomy or unknown menopausal status (n=38). There were no outlying values in these  
216 EPIC controls. In sum, our analysis included a total of 473 premenopausal EPIC controls.

217

## 218 **CA125 measurements**

219 In NEC controls we measured CA125 using the CA125II radioimmunoassay (Centocor,  
220 Malvern, PA) at the CER Lab at Boston Children's Hospitals. We assessed the reproducibility of  
221 the assay by including five blinded aliquots of a uniform quality control pool in each of the 46  
222 assay batches. The average of the coefficients of variation (CV) was 1%. In EPIC controls and in  
223 a subset of NEC controls, we previously measured CA125 using the volume-effective highly  
224 sensitive multiplex platform (Meso Scale Discovery (MSD), Gaithersburg, MD) in the Genital  
225 Tract Biology Laboratory at Brigham and Women's Hospital(17). The average CV across the  
226 assay batches was 19%.

227



228 **Candidate predictors**

229 We selected factors that have been previously reported to be associated with CA125 in at least one  
230 prior study(7-9,11), ovarian cancer risk factors(18), as well as several factors which were  
231 biologically plausible to be associated with CA125(10). Those included age at blood draw, race,  
232 body mass index (BMI, kg/m<sup>2</sup>), smoking status (never, former, current), pack-years calculated by  
233 number of packs of cigarettes per day multiplied by the number of years a person had smoked, age  
234 at menarche, oral contraceptive use and its duration (months), parity, self-reported endometriosis,  
235 tubal ligation, family history of ovarian cancer, prior personal cancer diagnosis, caffeine intake  
236 (mg), genital powder use, infertility, number of miscarriages, ectopic pregnancy, ever use of  
237 intrauterine device (IUD), fibroids, menstrual cycle regularity and days between last menstrual  
238 period (LMP) and blood draw(7-11). Furthermore, we evaluated additional variables related to  
239 menstrual characteristics and pregnancy timing: cycle length, days with menstrual bleeding,  
240 dysmenorrhea, age at first live birth, age at last live birth, and years since last live birth.

241

242 **Statistical analyses**

243 We log-transformed CA125 values to achieve a normal distribution. With this transformation, the  
244 distribution of log-transformed CA125 was normally distributed with skewness of 0.35 and  
245 kurtosis of 0.34, with a bell-shaped histogram. .

246

247 *Recalibration of CA125*

248 Since the EPIC samples had CA125 measured using an alternate assay (MSD assay) with a  
249 different scale, we used recalibration to rescale these measurement results to be comparable to the  
250 CA125II assay values. We recalibrated the EPIC CA125 values based on 187 NEC premenopausal  
251 controls with CA125 measurements on both CA125II and MSD assays using the drift correction  
252 method(19). We regressed the log-transformed MSD assay values to the log-transformed CA125II  
253 assay values and used the intercept and effect estimates of the model to calculate the recalibrated  
254 CA125II assay values based on the measured MSD assay values for all premenopausal controls in  
255 EPIC and used the recalibrated values in our analyses.

256 *Predictors of CA125 in premenopausal women*

257 First, we evaluated the association between individual candidate predictors and CA125 using linear  
258 or logistic regression adjusted for continuous age. We used effect estimates of the linear regression  
259 for each predictor to calculate the percent change in CA125 levels, calculated as  $[\exp(\beta) - 1]$   
260  $\times 100$  for a 1-unit change in the predictor. We determined the optimal modeling of continuous  
261 variables (age, BMI, age at menarche, duration of OC use, parity, and smoking pack-years) using  
262 restricted cubic splines to test for linearity(20). We used categorical variables for age, dichotomous  
263 variable for age at menarche, and piecewise linear spline with single knot for BMI since these  
264 variables were non-linearly associated with log-transformed CA125. We created composite  
265 categorical variables for OC use and duration and smoking status and pack-years, and compared  
266 nested models using likelihood ratio test and non-nested models using the Akaike information  
267 criterion and Vuong test(21). Based on these evaluations, candidate predictors were modeled as  
268 follows: age at blood draw (categorical, by 10 year intervals from age 30), race (white, non-white),  
269 BMI (piecewise linear spline model with single knot at 27), height (continuous, centered at 165),  
270 smoking status(categorical, never, former, current) and pack-years (continuous, never smokers,  
271 pack-years among former smokers, pack-years among current smokers), age at menarche (age 12  
272 and under, above 12), duration of OC use (continuous, including never users), parity (continuous),  
273 endometriosis (no, yes), tubal ligation (no, yes), family history of ovarian cancer (no, yes), prior  
274 personal cancer diagnosis (no, yes), caffeine intake (quartiles), genital powder use (no use, body  
275 use, genital use), infertility (no, yes), number of miscarriages (0, 1, 2, 3 or more), ectopic  
276 pregnancy (no, yes), intrauterine device use (never, ever), fibroids (no, yes), menstrual cycle  
277 regularity (regular, irregular) and predicted phase of the menstrual cycle (early follicular, late  
278 follicular, peri-ovulatory, luteal, long cycle, irregular, missing) based on the number of days  
279 between the last menstrual period and blood draw.

280

281 *Prediction modeling*

282 Overall, we developed CA125 prediction models (linear and dichotomous) in NEC and conducted  
283 external validation in EPIC (Figure 1). We used cross-validation to conduct internal validation of  
284 the model developed in NEC. Since information on some of the predictors were not collected in

285 EPIC, we developed an abridged model restricted to variables available in EPIC from the final  
286 model, and then validated the abridged model in EPIC.

287

### 288 *Linear model*

289 First, we developed a linear CA125 prediction model of log-transformed CA125 in NEC. We used  
290 stepwise linear regression analysis using  $p < 0.15$  as significance level for entry and retention in the  
291 model. In our primary prediction modeling, we used missing indicators for menstrual phase at  
292 blood draw due to a proportion of missing values. For variables with a limited number of missing  
293 values (age at menarche ( $n=1$ ), caffeine intake ( $n=23$ ), menstrual cycle length ( $n=23$ )), women  
294 with missing values were excluded. Age was forced in the model and the r-squared was calculated  
295 for the final prediction model, adjusted for study phase (1992-1997, 1998-2003, 2003-2008) and  
296 center (Massachusetts, New Hampshire). In addition, we calculated a delta r-squared that excluded  
297 the variability explained by study phase and center as these were matching factors in NEC and  
298 were forced into the model(22). The predicted log-transformed CA125 values in NEC were  
299 calculated using the effect estimates from the final prediction model. We evaluated the  
300 performance of the model by calculating the Pearson correlation coefficient to assess how well the  
301 predicted and the observed CA125 values agreed (i.e. calibration). We used 5-fold cross-validation  
302 to assess for overfitting in NEC and calculated the average r-squared across all sampled  
303 datasets(23). To evaluate potential bias due to missing data of candidate predictors, we conducted  
304 a sensitivity analysis restricted to women who had no missing predictors. We also conducted  
305 multiple imputation by chained equations (MICE) to impute the missing variables conditional on  
306 all of the predictors and outcome(24). We allowed 100 iterations and generated 20 imputed  
307 datasets. We applied the final prediction model in the 20 imputed datasets using the methods  
308 described and pooled the results of the model estimates using the Rubin's rules(25).

309

310 For external validation, we sought to validate our linear CA125 prediction model in EPIC.  
311 However, some of our key predictors (endometriosis and fibroids) were not collected in EPIC or  
312 were missing in majority of women (tubal ligation), thus, we validated an abridged model  
313 restricted to variables available in EPIC. First, among the predictors selected in the final model  
314 developed in NEC, we identified predictors available both in NEC and EPIC. We next ran the  
315 abridged model in NEC restricting to those variables available in both NEC and EPIC. We used

316 the effect estimates from this model to calculate the predicted value of log-transformed CA125 in  
317 the EPIC samples. We calculated the Pearson correlation coefficient between the predicted and the  
318 observed log-transformed CA125 to assess agreement and compare to that in the discovery dataset.  
319 We plotted the predicted versus the observed log-transformed CA125 for visual assessment.

320

### 321 *Dichotomous model*

322 Next, we developed and validated a dichotomous prediction model of elevated CA125 defined as  
323 having CA125  $\geq 35$  U/ mL following the same method used for developing the linear CA125  
324 prediction model described above but using logistic stepwise regression analysis. We evaluated  
325 the performance of the model by calculating the area under the curve (AUC) in the NEC  
326 (discovery) and EPIC (validation).

327

328 All statistical analyses were performed using SAS version 9.4 (SAS Institute Inc., Cary, NC),  
329 STATA version 12.1 (StataCorp, College Station, TX), and R version 3.4.3.

330

## 331 **Results**

332 Overall, the mean CA125 values were 17.3 U/ mL in 815 NEC controls and 14.9 U/ mL in 473  
333 EPIC controls after recalibration. The baseline characteristics of NEC and EPIC premenopausal  
334 women were similar except age at blood draw, race, age at menarche, OC use, current hormone  
335 use, infertility, parity, and tubal ligation were significantly different. (Supplementary Table S1).

336

### 337 *Recalibration of CA125*

338 We recalibrated the CA125 values in EPIC using the recalibration model based on 187 NEC  
339 premenopausal controls with both CA125II and MSD assay measurements. These two  
340 measurements were highly correlated ( $r=0.96$ , 95%CI 0.94, 0.97). After recalibration, the  
341 measured and recalibrated CA125 values also showed high correlation ( $r=0.95$ , 95%CI 0.93, 0.96).  
342 The recalibration model showed high performance in general.

343

344 *Predictors of CA125 in premenopausal women*

345 Age at blood draw was non-linearly associated with CA125, with women younger than 30 or more  
346 than 50 years old having significantly lower CA125 than those aged 30-39 years (Table 1). In age-  
347 adjusted models, menstrual phase at blood draw was significantly associated with CA125 levels,  
348 with early follicular phase levels being 8 to 21% higher than in other menstrual phases.  
349 Endometriosis and fibroids were associated with significantly higher CA125 levels, with 21 and  
350 13% difference, respectively, compared to those who did not have the condition. Current hormonal  
351 contraception use and tubal ligation were significantly associated with lower CA125 levels, with  
352 -16% and -11% difference, respectively. Cycle length, days with menstrual bleeding,  
353 dysmenorrhea, age at first live birth, age at last live birth, and years since last live birth were not  
354 significantly associated with CA125 levels in premenopausal women. Similar predictors were  
355 significantly associated with CA125 levels in the dichotomous model (Supplementary Table S2).

356

357 *Linear CA125 prediction modeling*

358 The final linear CA125 prediction model included age at blood draw, race, tubal ligation,  
359 endometriosis, menstrual phase at blood draw, and fibroids, with an r-squared of 0.07 (95%CI  
360 0.02, 0.09) (Table 2). The association between individual predictors and CA125 were similar in  
361 univariate and multivariate adjusted models. The r-squared of this full linear model when  
362 conducting 5-fold cross-validation was 0.02. When we restricted the analysis to the 498 controls  
363 with complete information on all predictors and applied the final linear CA125 prediction model,  
364 the r-squared was 0.12 (95%CI 0.05, 0.15) (Supplementary Table S3). When restricting to women  
365 with complete information on all predictors, the same variables were retained in the final model.  
366 For all the models, the delta r-squared, which subtracts the variance attributable to study phase and  
367 center from the total variance, was similar to the r-squared reported above. When evaluating the  
368 final continuous model in multiple imputed datasets in NEC, the beta coefficients, standard errors,  
369 and the r-squared were similar to the original model (Supplementary Table S3). The small  
370 differences in the measures of association when running the final model in the dataset using  
371 missing indicators, dataset restricted to those with complete information on all potential predictors,  
372 and multiple imputed datasets suggest that the missingness of menstrual phase at blood draw do

373 not largely influence the results. We also observed similar performance of the model when  
374 including all significant predictors in the univariate analyses, suggesting that the final model  
375 included important key predictors. Predicted log-transformed CA125 calculated based on the final  
376 model and the observed log-transformed CA125 were weakly correlated with a Pearson correlation  
377 coefficient of 0.26 (95%CI 0.19, 0.33) (Figure 2A).

378

379 For external validation, we developed an abridged linear CA125 prediction model which included  
380 age at blood draw, race, and menstrual phase at blood draw with r-squared of 0.05 (95%CI 0.01,  
381 0.07) in NEC (Table 2). Using the measures of association from this abridged model, we calculated  
382 the predicted log-transformed CA125 values in EPIC. The predicted log-transformed CA125  
383 values had a similar correlation with the observed log-transformed CA125 values in EPIC ( $r=0.22$   
384 (95%CI 0.13, 0.31)) as in the NEC abridged linear model ( $r=0.22$  (95%CI 0.15, 0.29)) (Figure 2B,  
385 2C). The spread of the predicted CA125 values in Figure 2 are much smaller than the spread of  
386 the observed CA125 values because the linear prediction model only explains a small proportion  
387 of the total variance of the observed CA125 values.

388

### 389 *Dichotomous CA125 prediction modeling*

390 The final dichotomous prediction model to predict women with  $CA125 \geq 35$  U/ mL included age  
391 at blood draw, tubal ligation, endometriosis, prior personal cancer diagnosis, family history of  
392 ovarian cancer, number of miscarriages, menstrual phase at blood draw, and smoking status and  
393 duration with an AUC of 0.83 (95%CI 0.77, 0.89) (Table 3, Figure 3). For menstrual phase at  
394 blood draw, we collapsed the other phase and irregular menstruation categories because few  
395 individuals had  $CA125 \geq 35$  U/mL in these groups. The association between individual predictors  
396 and CA125 were similar in univariate and multivariate adjusted models. The AUC of this full  
397 dichotomous model when conducting 5-fold cross-validation was 0.67. When we restricted the  
398 analysis to the 498 controls with complete information on all predictors and applied the final  
399 dichotomous model, the AUC was 0.84 (95%CI 0.76, 0.93) (Supplementary Table S4). When we  
400 conducted variable selection process using stepwise regression among women with complete  
401 information on all predictors, similar predictors were retained except number of miscarriages and

402 smoking status, resulting with an AUC of 0.79 (95%CI 0.69, 0.89). When evaluating the model in  
403 the multiple imputed datasets in NEC, the odds ratios and the AUC were largely similar to the  
404 primary analysis (Supplementary Table S4). We also observed a similar performance of the model  
405 when including all significant predictors from the univariate analyses, suggesting that the final  
406 model included important key predictors. We also considered using 65 U/mL cutoff which has  
407 been proposed for premenopausal women(26), but were limited with five controls who had CA125  
408 greater than 65 U/mL so were not able to investigate further.

409  
410 For external validation, we developed an abridged model, which included age at blood draw,  
411 number of miscarriages, menstrual phase (collapsing those on hormones, blood draw at other  
412 phase, and having irregular menstruation due to power), and smoking status with an AUC of 0.73  
413 (95%CI 0.65, 0.81) in NEC (Table 3, Figure3). When we applied this model to EPIC using  
414 recalibrated CA125 value of 35 U/ mL as cutoff, the AUC was 0.78 (95%CI 0.67, 0.89) (Figure  
415 3).

416

## 417 **Discussion**

418 This is the largest population-based study to develop and validate CA125 prediction models among  
419 healthy premenopausal women considering both continuous levels as well as those over current  
420 clinical threshold of 35 U/ mL. Although, the model predicting continuous CA125 only explained  
421 a small percent of the total variability, the model did show comparable correlations between  
422 predicted and observed levels in EPIC, suggesting the validity of the model. Conversely, the AUC  
423 for predicting elevated CA125 ( $\geq 35$  U/ mL) was relatively high in NEC and validated in EPIC.

424

425 Age was non-linearly associated with CA125 in our study, which is consistent with our prior study  
426 in EPIC in which we observed an inverse U-shaped association between age and CA125 levels  
427 among premenopausal women(9). Similarly, non-white race was associated with significantly  
428 lower CA125, which was consistent with prior studies in postmenopausal women(7,8), suggesting  
429 the need for different thresholds for minority populations. Unfortunately, we were underpowered

430 to evaluate differences in prediction models between racial subgroups, though others have  
431 described differences in CA125 levels between Black and Asian women(7,8).

432

433 Factors related to menstruation were strongly related to CA125. Specifically, an early follicular  
434 phase blood draw was significantly associated with higher CA125 levels and strong predictor of  
435 CA125 in our final model, which was consistent with previous reports(27). This association is  
436 likely driven by MUC16 expression on the endometrium and endometrial shedding during early  
437 follicular phase which may lead to higher circulating CA125 levels(10). This could explain the  
438 increased CA125 levels in women with fibroids, since fibroids are known to increase menstrual  
439 bleeding(28). In contrast, MUC16 expression on the endometrium may explain lower CA125  
440 levels among women with a tubal ligation as this procedure would prevent retrograde menstruation,  
441 which occurs in approximately 85% of women during menstruation(29), leading to systemic  
442 exposure to the antigen. Factors related to infertility, particularly endometriosis, were also related  
443 to substantially higher CA125 levels, consistent with prior studies(30,31). A similar mechanism is  
444 likely responsible as endometriosis leads to ectopic endometrial tissue usually in the peritoneal  
445 cavity.

446

447 Our linear CA125 prediction model explained 7% of the variability in CA125 but showed  
448 moderate validation in EPIC, whereas our dichotomous CA125 prediction model had better  
449 predictive ability with good validation. These results suggest that the variability of CA125 may be  
450 small in general but change dramatically by certain factors such as menstrual phase and  
451 endometriosis, and therefore the dichotomous prediction model performed better. We decided to  
452 use a standard log-linear model for developing the linear CA125 prediction model because the  
453 distribution of log-transformed CA125 was normally distributed with low kurtosis and skewness.  
454 When we included all significant predictors in the univariate analyses, both linear and dichotomous  
455 models showed similar performance compared to our final model having fewer predictors,  
456 suggesting some predictors were correlated.

457



458 Interestingly, some factors, such as fibroids and race were only significantly associated with  
459 continuous CA125 and some factors, such as prior personal cancer diagnosis and family history of  
460 ovarian cancer were only significantly associated with elevated CA125 (above 35 U/ mL). We  
461 suspect more predictors were selected in the final dichotomous CA125 prediction model because  
462 the association between exposures and CA125 were non-linear.

463

464 The major strength of our study is that we had two large independent population-based studies  
465 with detailed information on candidate predictors of CA125 to develop and validate CA125  
466 prediction models among premenopausal women. However, there are several limitations to our  
467 study. First, we had missing data on several variables. While we used missing indicators for our  
468 main analysis, our sensitivity analyses restricting to those with complete information on all  
469 predictors and using multiple imputation showed similar results, suggesting that the method for  
470 handling missing data did not influence overall results. In addition, we evaluated the performance  
471 of our prediction models using cross-validation and conducting external validation in an  
472 independent dataset, in which all the results were similar, suggesting a parsimonious model.  
473 Secondly, we were not able to validate the full prediction models in the independent dataset.  
474 Although we were only able to validate an abridged model in EPIC lacking tubal ligation and  
475 endometriosis, we expect the model performance to be better and closer to what we would have  
476 observed in NEC if we had information on all predictors. Thirdly, our model could be missing  
477 unknown predictors of CA125 since we restricted the candidate predictors to those previously  
478 described, which were mostly conducted among postmenopausal women. The relatively low r-  
479 squared of the final linear CA125 prediction model suggest that other candidate predictors may  
480 exist, such as genetic factors, common medications, and dietary factors, opening new opportunities  
481 for future studies. While hysterectomy has been previously described as a predictor of CA125,  
482 only few participants in NEC had hysterectomy. Given their ambiguous menopausal status we  
483 excluded them from current analysis of premenopausal women. Lastly, the model performance in  
484 EPIC could be underestimated because NEC and EPIC used different assays to measure CA125.  
485 However, the CA125 values of the two assays were highly correlated ( $r=0.96$ ) and the predicted  
486 CA125 values calculated using the recalibration model were also very highly correlated with the  
487 observed CA125II assay values ( $r=0.95$ ).

488

489 In summary, we developed and validated CA125 prediction models among premenopausal women  
490 in two independent studies that further our understanding of factors that influence CA125 levels  
491 and can therefore be used to optimize ovarian cancer screening with CA125. While performance  
492 of population-level screening for ovarian cancer in premenopausal women may be limited due to  
493 the lower incidence of ovarian cancer in this age range, approximately 30% of ovarian cancers are  
494 diagnosed before age 55. Furthermore, the impact of ovarian cancer in younger women results in  
495 potentially greater social, emotional, and economic impact. Further studies are needed to identify  
496 new predictors of CA125 to improve the model and to understand the predictors of changes in  
497 CA125 over time based on personal characteristics.

498

#### 499 **Acknowledgements**

500 The authors would like to acknowledge the participants and staff of the New England Case Control  
501 Study and the European Prospective Investigation into Cancer and Nutrition study.

502

#### 503 **References**

- 504 1. Howlader N NA, Krapcho M, Miller D, Bishop K, Kosary CL, Yu M, Ruhl J, Tatalovich  
505 Z, Mariotto A, Lewis DR, Chen HS, Feuer EJ, Cronin KA (eds). SEER Cancer Statistics  
506 Review, 1975-2014. National Cancer Institute. Bethesda, MD, 2017.
- 507 2. Torre LA, Islami F, Siegel RL, Ward EM, Jemal A. Global Cancer in Women: Burden  
508 and Trends. *Cancer Epidemiol Biomarkers Prev.* 2017;26(4):444-457.
- 509 3. Duffy MJ, Bonfrer JM, Kulpa J, Rustin GJ, Soletormos G, Torre GC, et al. CA125 in  
510 ovarian cancer: European Group on Tumor Markers guidelines for clinical use. *Int J*  
511 *Gynecol Cancer.* 2005;15(5):679-691.
- 512 4. Bast RC, Jr., Klug TL, St John E, Jenison E, Niloff JM, Lazarus H, et al. A  
513 radioimmunoassay using a monoclonal antibody to monitor the course of epithelial  
514 ovarian cancer. *N Engl J Med.* 1983;309(15):883-887.
- 515 5. Buys SS, Partridge E, Black A, Johnson CC, Lamerato L, Isaacs C, et al. Effect of  
516 screening on ovarian cancer mortality: the Prostate, Lung, Colorectal and Ovarian  
517 (PLCO) Cancer Screening Randomized Controlled Trial. *JAMA.* 2011;305(22):2295-  
518 2303.
- 519 6. Jacobs IJ, Menon U, Ryan A, Gentry-Maharaj A, Burnell M, Kalsi JK, et al. Ovarian  
520 cancer screening and mortality in the UK Collaborative Trial of Ovarian Cancer

521 Screening (UKCTOCS): a randomised controlled trial. *Lancet*. 2016;387(10022):945-  
522 956.

523 7. Pauler DK, Menon U, McIntosh M, Symecko HL, Skates SJ, Jacobs IJ. Factors  
524 influencing serum CA125II levels in healthy postmenopausal women. *Cancer Epidemiol*  
525 *Biomarkers Prev*. 2001;10(5):489-493.

526 8. Johnson CC, Kessel B, Riley TL, Ragard LR, Williams CR, Xu JL, et al. The  
527 epidemiology of CA-125 in women without evidence of ovarian cancer in the Prostate,  
528 Lung, Colorectal and Ovarian Cancer (PLCO) Screening Trial. *Gynecol Oncol*.  
529 2008;110(3):383-389.

530 9. Fortner RT, Vitonis AF, Schock H, Husing A, Johnson T, Fichorova RN, et al. Correlates  
531 of circulating ovarian cancer early detection markers and their contribution to  
532 discrimination of early detection models: results from the EPIC cohort. *J Ovarian Res*.  
533 2017;10(1):20.

534 10. Haridas D, Ponnusamy MP, Chugh S, Lakshmanan I, Seshacharyulu P, Batra SK.  
535 MUC16: molecular analysis and its functional implications in benign and malignant  
536 conditions. *FASEB J*. 2014;28(10):4183-4199.

537 11. Westhoff C, Gollub E, Patel J, Rivera H, Bast R, Jr. CA 125 levels in menopausal  
538 women. *Obstet Gynecol*. 1990;76(3 Pt 1):428-431.

539 12. Vitonis AF, Titus-Ernstoff L, Cramer DW. Assessing ovarian cancer risk when  
540 considering elective oophorectomy at the time of hysterectomy. *Obstet Gynecol*.  
541 2011;117(5):1042-1050.

542 13. Willett WC, Sampson L, Stampfer MJ, Rosner B, Bain C, Witschi J, et al.  
543 Reproducibility and validity of a semiquantitative food frequency questionnaire. *Am J*  
544 *Epidemiol*. 1985;122(1):51-65.

545 14. Rimm EB, Giovannucci EL, Stampfer MJ, Colditz GA, Litin LB, Willett WC.  
546 Reproducibility and validity of an expanded self-administered semiquantitative food  
547 frequency questionnaire among male health professionals. *Am J Epidemiol*.  
548 1992;135(10):1114-1126; discussion 1127-1136.

549 15. Rosner B. Percentage points for a generalized ESD many-outlier procedure. .  
550 *Technometrics* 1983;25:165-172.

551 16. Riboli E, Hunt KJ, Slimani N, Ferrari P, Norat T, Fahey M, et al. European Prospective  
552 Investigation into Cancer and Nutrition (EPIC): study populations and data collection.  
553 *Public Health Nutr*. 2002;5(6B):1113-1124.

554 17. Terry KL, Schock H, Fortner RT, Husing A, Fichorova RN, Yamamoto HS, et al. A  
555 Prospective Evaluation of Early Detection Biomarkers for Ovarian Cancer in the  
556 European EPIC Cohort. *Clin Cancer Res*. 2016;22(18):4664-4675.

557 18. Wentzensen N, Poole EM, Trabert B, White E, Arslan AA, Patel AV, et al. Ovarian  
558 Cancer Risk Factors by Histologic Subtype: An Analysis From the Ovarian Cancer  
559 Cohort Consortium. *J Clin Oncol*. 2016;34(24):2888-2898.

560 19. Eliassen AH, Hendrickson SJ, Brinton LA, Buring JE, Campos H, Dai Q, et al.  
561 Circulating carotenoids and risk of breast cancer: pooled analysis of eight prospective  
562 studies. *J Natl Cancer Inst*. 2012;104(24):1905-1916.

563 20. Durrleman S, Simon R. Flexible regression models with cubic splines. *Stat Med*.  
564 1989;8(5):551-561.

565 21. Clarke KA. Testing Nonnested Models of International Relations: Reevaluating Realism.  
566 *American Journal of Political Science*. 2001;45(3):724-744.

- 567 22. Rosner B. *Fundamentals of Biostatistics*. Eighth Edition ed: CENGAGE Learning; 2016.
- 568 23. Tworoger SS, Zhang X, Eliassen AH, Qian J, Colditz GA, Willett WC, et al. Inclusion of  
569 endogenous hormone levels in risk prediction models of postmenopausal breast cancer. *J*  
570 *Clin Oncol*. 2014;32(28):3111-3117.
- 571 24. van Buuren S. Multiple imputation of discrete and continuous data by fully conditional  
572 specification. *Stat Methods Med Res*. 2007;16(3):219-242.
- 573 25. DB R. *Multiple Imputation for Nonresponse in Surveys*.: John Wiley and Sons: New  
574 York; 1987.
- 575 26. Eltabbakh GH, Belinson JL, Kennedy AW, Gupta M, Webster K, Blumenson LE. Serum  
576 CA-125 measurements > 65 U/mL. Clinical value. *J Reprod Med*. 1997;42(10):617-624.
- 577 27. Bon GG, Kenemans P, Dekker JJ, Hompes PG, Verstraeten RA, van Kamp GJ, et al.  
578 Fluctuations in CA 125 and CA 15-3 serum concentrations during spontaneous ovulatory  
579 cycles. *Hum Reprod*. 1999;14(2):566-570.
- 580 28. Bischof P, Galfetti MA, Seydoux J, von Hostenpenthal JU, Campana A. Peripheral CA 125  
581 levels in patients with uterine fibroids. *Hum Reprod*. 1992;7(1):35-38.
- 582 29. D'Hooghe TM, Bambra CS, Raeymaekers BM, Koninckx PR. Increased prevalence and  
583 recurrence of retrograde menstruation in baboons with spontaneous endometriosis. *Hum*  
584 *Reprod*. 1996;11(9):2022-2025.
- 585 30. Muyldermans M, Cornillie FJ, Koninckx PR. CA125 and endometriosis. *Hum Reprod*  
586 *Update*. 1995;1(2):173-187.
- 587 31. Mol BW, Bayram N, Lijmer JG, Wiegerinck MA, Bongers MY, van der Veen F, et al.  
588 The performance of CA-125 measurement in the detection of endometriosis: a meta-  
589 analysis. *Fertil Steril*. 1998;70(6):1101-1108.

590

591

592

593

594

595

596

597

598

599

**Table 1. Association between predictors and CA125 levels in premenopausal women without ovarian cancer in the New England Case Control Study (n=815)**

Variables	N (%)	Mean CA125 (95%CI) <sup>a</sup>	Differences in CA125 levels <sup>b</sup>	p-value <sup>b</sup>
<b>Age, years</b>				
< 30	69 (8)	13 (11,14)	-20%	0.001
30- 39	215 (26)	16 (15, 17)	ref	
40- 49	374 (46)	15 (15, 16)	-5%	0.20
50+	157 (19)	15 (13, 16)	-10%	0.05
P-trend				0.65
<b>Race</b>				
white	778 (95)	15 (15, 16)	ref	
non-white	37 (5)	14 (12, 16)	-10%	0.20
<b>BMI, kg/m<sup>2</sup></b>				
< 20	76 (9)	15 (14, 17)	1%	0.90
20- < 25	414 (51)	15 (14, 16)	ref	
25- < 30	206 (25)	15 (14, 16)	-1%	0.77
30- < 35	84 (10)	16 (14, 17)	3%	0.6
35+	35 (4)	17 (14, 20)	13%	0.16
P-trend				0.30
<b>Height, cm</b>				
< 160	182 (22)	16 (14, 17)	ref	
160- < 165	218 (27)	14 (13, 15)	-8%	0.08
165- < 170	231 (28)	16 (15, 17)	1%	0.85
170- < 175	128 (16)	15 (14, 17)	-1%	0.84
175+	56 (7)	14 (12, 16)	-8%	0.25
P-trend				0.88
<b>Smoking</b>				
never	421 (52)	15 (14, 16)	ref	
former	269 (33)	16 (15, 17)	5%	0.19
current	125 (15)	14 (13, 15)	-7%	0.16
<b>Smoking status and duration</b>				
never smokers	421 (52)	15 (14, 16)	ref	
< 5 packyears among former	140 (17)	16 (15, 17)	5%	0.30
5- < 15 packyears among former	79 (10)	16 (15, 18)	8%	0.20
15+ packyears among former	50 (6)	15 (13, 17)	1%	0.92
< 5 packyears among current	23 (3)	15 (12, 19)	1%	0.93
5- < 15 packyears among current	46 (6)	14 (12, 16)	-6%	0.47
15+ packyears among current	56 (7)	13 (12, 15)	-11%	0.11
<b>Age at menarche, years</b>				

≤ 12	386 (47)	16 (15, 16)	ref	
13+	428 (53)	15 (14, 15)	-6%	0.09
<b>Oral contraceptive use</b>				
Never	180 (22)	15 (14, 16)	ref	
Ever	635 (78)	15 (15, 16)	5%	0.25
<b>Duration of oral contraceptive use among ever users, years</b>				
< 2	140 (22)	17 (16, 19)	ref	
2 - 3	133 (21)	15 (14, 17)	-10%	0.07
4 - 5	103 (16)	15 (14, 17)	-11%	0.06
6 - 9	125 (20)	15 (13, 16)	-14%	0.01
10+	134 (21)	14 (13, 16)	-16%	0.005
P-trend				0.01
<b>Current hormonal contraception use<sup>c</sup></b>				
no	427 (82)	15 (15, 16)	ref	
yes	94 (18)	13 (12, 14)	-16%	0.003
<b>Menstrual phase at time of blood draw (days)</b>				
early follicular (0 - 7)	126 (15)	17 (15, 18)	ref	
late follicular (8 - 11)	56 (7)	16 (14, 18)	-8%	0.31
peri-ovulatory (12 - 16)	73 (9)	15 (13, 16)	-13%	0.06
luteal (17 - 35)	155 (19)	14 (13, 16)	-14%	0.01
long cycle (36+)	63 (8)	13 (12, 15)	-21%	0.002
irregular	65 (8)	13 (12, 15)	-21%	0.002
missing	277 (34)	16 (15, 17)	-8%	0.15
<b>Cause of infertility</b>				
none	649 (80)	15 (15, 16)	ref	
male factor	20 (2)	17 (14, 21)	12%	0.32
tubal	13 (2)	15 (11, 20)	-1%	0.92
endometriosis	13 (2)	23 (17, 30)	50%	0.004
ovulatory	14 (2)	14 (11, 18)	-7%	0.59
unknown	106 (13)	14 (13, 15)	-9%	0.08
<b>Number of miscarriages</b>				
0	625 (77)	15 (15, 16)	ref	
1	137 (17)	15 (13, 16)	-5%	0.29
2+	53 (7)	14 (13, 17)	-6%	0.41
<b>Ectopic pregnancy</b>				
no	795 (98)	15 (15, 16)	ref	
yes	20 (2)	16 (13, 20)	5%	0.67
<b>Parity</b>				
0	214 (26)	15 (14, 16)	ref	
1	140 (17)	15 (14, 16)	1%	0.82
2	264 (32)	15 (14, 16)	3%	0.53

3+	197 (24)	15 (14, 16)	2%	0.66
<b>Tubal ligation</b>				
no	674 (83)	15 (15, 16)	ref	
yes	141 (17)	14 (13, 15)	-11%	0.01
<b>Intrauterine device use</b>				
never	689 (85)	15 (14, 16)	ref	
ever	126 (15)	16 (15, 17)	6%	0.26
<b>Unilateral oophorectomy</b>				
no	808 (99)	15 (15, 16)	ref	
yes	7 (1)	19 (13, 27)	23%	0.28
<b>Endometriosis</b>				
no	764 (94)	15 (14, 15)	ref	
yes	51 (6)	18 (16, 21)	21%	0.01
<b>Fibroids</b>				
no	737 (90)	15 (14,16)	ref	
yes	78 (10)	17 (15,19)	13%	0.04
<b>Prior personal cancer diagnosis</b>				
no	778(95)	15 (15, 16)	ref	
yes	37(4)	17 (15, 20)	15%	0.10
<b>Family history of ovarian cancer</b>				
no	795 (98)	15 (15, 16)	ref	
yes	20 (2)	17 (13, 21)	10%	0.40
<b>Genital powder use</b>				
no use	484 (59)	15 (15, 16)	ref	
body use	157 (19)	15 (14, 16)	-5%	0.30
genital use	174 (21)	15 (14, 16)	-3%	0.57
<b>Caffeine intake, mg</b>				
< 70.1	198 (25)	16 (14, 17)	ref	
70.1- < 169.6	198 (25)	14 (13, 15)	-9%	0.07
169.6- < 348.7	198 (25)	15 (14, 17)	0%	0.97
348.7+	198 (25)	16 (15, 17)	1%	0.90
P-trend				0.44

<sup>a</sup>geometric mean adjusted for age

<sup>b</sup>age-adjusted

<sup>c</sup>includes oral contraceptives and injections

601

602

603

604

**Table 2. Linear CA125 prediction model in premenopausal women using stepwise regression in the New England Case Control Study (NEC)**

Selected Predictors	Development of model in NEC (n=768) <sup>a b</sup>			Abridged model in NEC (n=768) <sup>a b</sup>		
	differences in CA125 levels	SE	p-value	differences in CA125 levels	SE	p-value
<b>Age, years</b>						
< 30	-18%	0.07	0.01	-19%	0.07	0.004
30- < 40	ref			ref		
40- < 50	-8%	0.04	0.06	-8%	0.04	0.05
50+	-10%	0.06	0.06	-11%	0.06	0.04
<b>Non-white race</b>	-13%	0.09	0.10	-14%	0.09	0.08
<b>Tubal ligation</b>	-10%	0.05	0.03			
<b>Endometriosis</b>	21%	0.07	0.01			
<b>Menstrual phase at time of blood draw</b>						
early follicular phase	ref			ref		
on hormonal contraceptives <sup>c</sup>	-28%	0.07	<.0001	-27%	0.07	<.0001
other cycle phase	-17%	0.06	0.002	-17%	0.06	0.001
irregular cycles	-30%	0.11	0.001	-30%	0.11	0.001
missing	-17%	0.10	0.05	-18%	0.10	0.04
<b>Fibroids</b>	13%	0.06	0.04			
<b>R-square<sup>d</sup></b>	0.07			0.05		

<sup>a</sup>adjusted for study phase and center

<sup>b</sup>exclude missing age at menarche(n=1), caffeine intake(n=23), and menstrual cycle length(n=23)

<sup>c</sup>includes oral contraceptives and injections

<sup>d</sup>delta r-square, subtracting the effect of center and study phase

605

606

607

608

609

610

611



**Table 3. Dichotomous CA125 prediction model of high CA125 ( $\geq 35$  U/ mL) in premenopausal women using stepwise regression in the New England Case Control Study (NEC)**

Selected predictors	Development of model in NEC(n=768) <sup>a b</sup>		Abridged model in NEC(n=768) <sup>a b</sup>	
	Odds ratio	95%CI	Odds ratio	95%CI
<b>Age, years</b>				
< 30	0.34	0.07, 1.61	0.29	0.06, 1.30
30- < 40	ref		ref	
40- < 50	0.34	0.15, 0.74	0.40	0.20, 0.82
50+	0.37	0.13, 1.04	0.44	0.17, 1.12
<b>Tubal ligation</b>	0.16	0.03, 0.73		
<b>Endometriosis</b>	3.08	1.09, 8.72		
<b>Prior personal cancer diagnosis</b>	5.41	1.70, 17.18		
<b>Family history of ovarian cancer</b>	13.32	3.42, 51.93		
<b>Number of miscarriages</b>				
0	ref		ref	
1	0.23	0.06, 0.89	0.36	0.11, 1.21
2+	0.18	0.02, 1.56	0.29	0.04, 2.18
<b>Menstrual phase at time of blood draw</b>				
regular menstruation + blood draw at early follicular phase on hormonal contraceptives <sup>c</sup>	ref		ref	
regular menstruation + blood draw at other phase /irregular menstruation	0.07	0.01, 0.63	0.35	0.15, 0.82
missing time at blood draw	0.49	0.19, 1.25	0.35	0.05, 2.45
<b>Smoking status</b>				
Never	ref		ref	
Former	2.15	0.97, 4.79	1.85	0.96, 3.57
Current	0.12	0.01, 1.51	0.47	0.13, 1.64
Pack-years among former smokers	0.99	0.94, 1.05		
Pack-years among current smokers	1.09	1.00, 1.18		
<b>AUC (95%CI)</b>	0.83	0.77, 0.89	0.73	0.65, 0.81

<sup>a</sup>adjusted for study phase and center

<sup>b</sup>exclude missing age at menarche(n=1), caffeine intake(n=23), and menstrual cycle length(n=23)

<sup>c</sup>includes oral contraceptives and injections

614 **Figure Legend**

615 **Figure 1. Study design of the development and validation of the CA125 prediction model**  
616 **using the New England Case Control Study (NEC) and European Prospective Investigation**  
617 **into Cancer and Nutrition Study (EPIC)**

618 We developed the CA125 prediction models (linear and dichotomous) using the New England  
619 Case Control Study (NEC) and conducted external validation using the European Prospective  
620 Investigation into Cancer and Nutrition Study (EPIC).

621

622 **Figure 2. Development and validation of linear CA125 prediction model in the New England**  
623 **Case Control Study (NEC) and European Prospective Investigation into Cancer and**  
624 **Nutrition Study (EPIC)**

625 The predicted versus the observed log-transformed CA125 values were plotted and Pearson  
626 correlation coefficient ( $r$ ) was calculated to assess the performance of the linear CA125 prediction  
627 model in the New England Case Control Study (NEC) and European Prospective Investigation  
628 into Cancer and Nutrition Study (EPIC). **A**, Linear CA125 prediction model performance in NEC.  
629 **B**, Abridged linear CA125 prediction model performance in NEC. **C**, Abridged linear CA125  
630 prediction model performance in EPIC.

631

632 **Figure 3. Development and validation of dichotomous CA125 prediction model the New**  
633 **England Case Control Study (NEC) and European Prospective Investigation into Cancer**  
634 **and Nutrition Study (EPIC)**

635 Receiver Operating Characteristic (ROC) curves were plotted and the Area Under the Curve  
636 (AUC) was calculated to assess the discriminatory performance of the dichotomous CA125  
637 prediction model in the New England Case Control Study (NEC) and European Prospective  
638 Investigation into Cancer and Nutrition Study (EPIC). Dichotomous CA125 prediction model  
639 performance in NEC (solid line), abridged dichotomous CA125 prediction model performance in

640 NEC (dashed line), abridged dichotomous CA125 prediction model performance in EPIC (dotted  
641 line).

642