

1 **External validation of ultrasound-based models for discrimination between benign and**
2 **malignant adnexal masses in Italy: the prospective multicenter IOTA phase 6 study**

3

4 Francesca Moro^{1,2}, Marina Momi², Valentina Bertoldo², Ashleigh Ledger³, Lasai Barreñada³, Jolien
5 Ceusters³, Davide Sturla⁴, Fabio Ghezzi⁴, Elisa Mor⁵, Letizia Fornari⁶, Antonella Vimercati⁷,
6 Saverio Tateo⁸, Marianna Roccio⁹, Rosalba Giacchello¹⁰, Roberta Granese¹¹, Daniela Garbin¹²,
7 Tiziana De Grandis¹³, Federica Piccini¹⁴, Patrizia Favaro¹⁵, Olga Petruccelli¹⁶, Anila Kardhashi¹⁷,
8 Ilaria Pezzani¹⁸, Patrizia Ragno¹⁹, Laura Falchi²⁰, Bruna Anna Virgilio²¹, Erika Fruscella²², Tiziana
9 Tagliaferri²³, Annibale Mazzocco²⁴, Floriana Mascilini², Francesca Ciccarone², Federica Pozzati²,
10 Wouter Froyman^{3,25}, Ben Van Calster³, Tom Bourne²⁶, Dirk Timmerman^{3,25}, Giovanni Scambia²,
11 *Lil Valentin^{27,28}, *Antonia Carla Testa²

12

13 ¹ UniCamillus, International Medical University, Rome, Italy
14 ² Fondazione Policlinico Universitario A. Gemelli IRCCS, Rome, Italy
15 ³ Department of Development and Regeneration, KU Leuven, Leuven, Belgium
16 ⁴ Ospedale Filippo del Ponte Ostetricia e Ginecologia, Varese, Lombardia, Italy
17 ⁵ Fondazione Poliambulanza Istituto Ospedaliero, Brescia, Lombardia, Italy
18 ⁶ Azienda Ospedaliero Universitaria Pisana, Pisa, Toscana, Italy
19 ⁷ Azienda Ospedaliero-Universitaria Consorziale Policlinico di Bari, Bari, Puglia, Italy
20 ⁸ Ospedale Santa Chiara di Trento, Trento, Trentino-Alto Adige, Italy
21 ⁹ Fondazione IRCCS Policlinico San Matteo, Pavia, Lombardia, Italy
22 ¹⁰ ASL CN1, Cuneo, Piemonte, Italy
23 ¹¹ Azienda Ospedaliera Universitaria Policlinico "G. MARTINO"
24 ¹² Ospedale di Santorso
25 ¹³ IRCCS Candiolo
26 ¹⁴ "Ramazzini" di Carpi
27 ¹⁵ Ospedale di Bussolengo
28 ¹⁶ Azienda Ospedaliera Universitaria O.O.R.R. Foggia
29 ¹⁷ Istituto Tumori, IRCCS, "Giovanni Paolo II"
30 ¹⁸ Presidio Ospedaliero di Treviso
31 ¹⁹ ASL AT - Asti
32 ²⁰ Azienda USL Toscana Centro
33 ²¹ Ospedale di Abano Terme
34 ²² Ospedale Santo Spirito - Roma
35 ²³ Ospedale San Maurizio di Bolzano
36 ²⁴ Presidio Ospedaliero di Montebelluna, Veneto, Italy
37 ²⁵ Department of Obstetrics and Gynecology, University Hospitals KU Leuven, Leuven, Belgium
38 ²⁶ Imperial College Healthcare NHS Trust, London, UK
39 ²⁷ Skåne University Hospital, Malmö, Sweden
40 ²⁸ Department of Clinical Sciences Malmö, Lund University, Sweden

41

42 *The last two authors contributed equally to the work.

43

44

45 **Corresponding authors:**

46 Francesca Moro

47 UniCamillus, International Medical University, Rome, Italy

48 Fondazione Policlinico Universitario A. Gemelli IRCCS, Rome, Italy

49 E-mail: morofrancy@gmail.com; francesca.moro@guest.policlinicogemelli.it;

50 Phone number: +39-3493905802

51

52 **Keywords:**

53 Ovarian neoplasm, Ultrasonography,

54

55

56

57

58

59

60

61

62

63

64

65

66

67

68 **Abstract**

69 **Objective** To prospectively validate the performance of the Risk of Malignancy Index (RMI),
70 International Ovarian Tumor Analysis (IOTA) Simple Rules Risk Model (SRRisk), IOTA
71 Assessment of Different NEoplasias in the adneXa (ADNEX) and the IOTA two-step strategy in
72 different types of ultrasound centers in Italy.

73 **Methods** This is a multicenter prospective observational study including regional referral centers
74 and district hospitals in Italy. Consecutive patients with an adnexal mass examined with ultrasound
75 by an IOTA certified ultrasound examiner with different levels of experience were included,
76 provided they underwent surgery < 180 days after the inclusion scan. Ultrasound examination was
77 performed transvaginally or transrectally and/or transabdominally based on the characteristics of the
78 women and masses. Reference standard was the histology of the adnexal mass following surgical
79 removal. Discrimination (area under receiver operating characteristic curve, AUROC), calibration,
80 and clinical utility were assessed to illustrate the diagnostic performance of the methods. The
81 performance of the models was also evaluated in predefined subgroups based on menopausal status,
82 type of center (oncology vs non-oncology) and ultrasound examiner's experience: [<500 scans
83 performed, 500-5000 scans performed, >5000 scans performed; European Federation of Societies
84 for Ultrasound in Medicine and Biology (EFSUMB) Level 1, Level 2, Level 3].

85 **Results** 1567 patients were recruited between May 2017 and March 2020 from 23 Italian centers.
86 After data cleaning and application of exclusion criteria, our study population consisted of 1431
87 patients in 21 Italian centers (10 oncological and 11 non-oncological). Based on histology,
88 995/1431 (69.5%) tumors were benign and 436/1431 (30.5%) were malignant (115/1431, 8.0%
89 borderline, 263/1431, 18.4% primary invasive, 58/1431, 4.1% metastatic tumors). For all IOTA
90 models (SRRisk, ADNEX with and without CA125, two step strategy with and without CA125),
91 the AUROC was between 0.91 (95% CI 0.88-0.93) and 0.92 (0.89-0.94). The AUROC was 0.85
92 (0.81-0.87) for RMI. The malignancy risk was slightly underestimated by all IOTA models, but

93 least so by SRRisk . All IOTA models had higher net benefit than RMI at risk thresholds from 1%
94 to 50%. AUROC was ≥ 0.90 for all IOTA models in all subgroups, while it ranged from 0.84 to 0.90
95 for RMI.

96 **Conclusions** SRRisk, ADNEX and the two step strategy with or without CA125 had similar and
97 good ability to distinguish benign from malignant adnexal tumours in patients examined by either
98 expert or non-expert ultrasound operators in Italy. Their discriminative performance and clinical
99 utility was superior to that of RMI.

100

101 **Introduction**

102 Ovarian cancer is the leading cause of death in women diagnosed with gynecological cancers.¹
103 Most ovarian cancers are diagnosed at an advanced stage and require treatment in high volume
104 centers by doctors with expertise in gynecological oncological surgery to optimize outcome.²⁻⁴
105 Correct preoperative characterization of adnexal masses is essential to decide on optimal
106 management: clinical and ultrasound follow-up, surgery in a local center, or referral to an oncology
107 center.^{5,6} Transvaginal ultrasound is the first line method for characterizing adnexal masses. If
108 performed by an expert, subjective assessment of the ultrasound images is the optimal method for
109 distinguishing benign from malignant masses.⁷⁻⁹ For less experienced ultrasound examiners there
110 are other methods. The Risk of Malignancy Index (RMI) is a scoring system using clinical and
111 ultrasound information that can be used to estimate the likelihood of an ovarian mass being
112 malignant.¹⁰ In some European countries RMI is widely used to triage women with an adnexal mass
113 for referral to an oncological center.¹¹⁻¹⁵ The International Ovarian Tumor Analysis (IOTA) group
114 has developed several ultrasound based methods that can be used to discriminate between benign
115 and malignant adnexal masses: the Benign Descriptors,^{16,17} the Simple Rules,¹⁸ four mathematical
116 models to calculate the individual risk of malignancy in an adnexal mass (logistic regression model
117 1, LR1, logistic regression model 2, LR2, the Simple Rules Risk model, SRRisk, and Assessment
118 of Different NEoplasias in the adneXa - ADNEX).¹⁹⁻²¹ ADNEX is a multinomial regression model
119 that calculates the probability of five outcome categories (benign, borderline, stage I primary
120 invasive ovarian malignancy, stage II-IV primary invasive ovarian malignancy, and metastasis from
121 another primary tumor).²¹ ADNEX can be used with or without CA125 as predictor.²¹ The IOTA
122 group now recommends the IOTA two-step strategy, which means that first the Benign Descriptors
123 are applied, and if these do not apply, ADNEX is used.¹⁷
124 The diagnostic performance of the IOTA methods and of RMI has been validated in prospective
125 and retrospective studies, but most validation studies tested the performance in the hands of

126 experienced ultrasound examiners.^{17,22–32} No prospective study included examiners with little
127 ultrasound experience and few included examiners with different levels of experience.^{33–35}
128 The primary aim of this study is to prospectively validate the diagnostic performance of RMI,
129 SRRisk, ADNEX and the IOTA two-step strategy in different types of ultrasound centers in Italy
130 both overall and in relevant subgroups. The secondary aims are to explore the multinomial
131 discrimination performance of ADNEX and the two step strategy, to validate the ability of the
132 Benign Descriptors to correctly classify an adnexal mass as benign, and to validate the classification
133 performance of the Simple Rules and of subjective assessment overall and in the subgroups based
134 on level of ultrasound experience.

135

136 **Methods**

137 *Study design and participants*

138 This is an Italian multicenter prospective external validation study of ultrasound based models to
139 discriminate between benign and malignant adnexal masses. The protocol was approved by the
140 Ethical Committee of the Fondazione Policlinico A. Gemelli, IRCCS (PROT 27665/16) and of each
141 participating center (Appendix 1). Written informed consent was obtained from all patients. The
142 study was conducted in accordance with the TRIPOD cluster guidelines.

143

144 *Patients*

145 Consecutive patients with a known or suspected adnexal mass examined with ultrasound by an
146 IOTA certified ultrasound examiner³⁶ and confirmed to have an adnexal mass judged not to be
147 physiological were eligible for inclusion provided they were expected to undergo surgical removal
148 of the mass. The patients were collected between May 2017 and March 2020. Exclusion criteria
149 were: patient's age <18 years, pregnant patients, patients with previous bilateral adnexectomy,
150 patients examined in centers that recruited < 10 patients, only transabdominal ultrasound

151 performed, surgery performed more than 180 days after the ultrasound examination, and denial or
152 withdrawal of informed consent.

153

154 *Data collection*

155 Information on age, parity, menopausal status and indication for the ultrasound examination was
156 prospectively collected, as well as information on type of hospital (private practice, local public
157 hospital, regional public hospital, or university hospital), type of center (oncological vs non-
158 oncological), and type of ultrasound center (general gynecologic outpatient clinic or specialized
159 ultrasound center). An oncological center was defined as a tertiary referral center with a dedicated
160 gynecological oncology unit. Information on the ultrasound system used, ultrasound examiner's
161 name and level of experience was also recorded. The level of ultrasound experience was based on
162 the number of gynecological scans in non-pregnant women that the examiner had performed at the
163 start of the study. Low experience was defined as <500 scans, intermediate experience as 500-5000
164 scans, and high experience as >5000 scans. We also recorded the level of experience according to
165 the European Federation of Societies for Ultrasound in Medicine and Biology (EFSUMB level 1, 2
166 or 3³⁷), and the number of ovarian masses that the operator had examined with ultrasound during
167 the preceding year (classified as <50, i.e. < 1 per week; 50-200, i.e. up to 4 per week; >200, i.e. > 4
168 per week).

169

170 *Ultrasound examination*

171 A standardized transvaginal (or transrectal if vaginal was not possible) ultrasound examination
172 including color or power Doppler ultrasound examination was performed, supplemented with
173 transabdominal ultrasound if transvaginal or transrectal ultrasound examination was not sufficient.
174 The IOTA examination and measurement technique were used, and the ultrasound findings were
175 described using the IOTA terminology.³⁸ Information on all the variables required for the IOTA
176 Benign Descriptors, Simple Rules, SRRisk, ADNEX and RMI were prospectively collected and

177 recorded. Results of subjective assessment were recorded as benign, borderline or malignant. The
178 degree of diagnostic confidence (certainly benign, probably benign, uncertain, probably malignant
179 or probably borderline, certainly malignant or certainly borderline) as well as the specific diagnosis
180 suggested by the ultrasound examiner and chosen from a list of pre-defined diagnoses were also
181 recorded.

182 If more than one adnexal mass was present, only the one with the most complex ultrasound
183 morphology was included in our statistical analysis. If the ultrasound morphology was similar in all
184 masses, the largest one or the one most easily accessible with ultrasound was used in our statistical
185 calculations. The management was decided by the referring clinician, who took into account clinical
186 symptoms, ultrasound results based on subjective evaluation of the ultrasound images (i.e. those
187 reported in the clinical ultrasound report), and results of other imaging modalities (e.g. computer
188 tomography or magnetic resonance imaging), tumor markers, and patient's preference.

189

190 *Reference standard*

191 Reference standard was the histology of the adnexal mass following surgical removal within 180
192 days after the ultrasound examination by laparotomy or laparoscopy as considered appropriate by
193 the surgeon. Borderline tumors were classified as malignant. The histology of the surgically
194 removed tumor was determined at each local center. Central pathology review was not performed,
195 because we found little differences between local and central pathology reports in a previous IOTA
196 study.¹⁹ Pathologists were blinded to ultrasound predictor variables and model predictions but might
197 have received information on the subjective assessment by the ultrasound examiner when clinically
198 relevant. The stage of malignant tumors was recorded using the classification of the International
199 Federation of Gynecology and Obstetrics (FIGO).³⁹

200

201 *Data cleaning*

202 Data collection was done through the web-based clinical data miner (CDM) software.⁴⁰ Patients
203 automatically received a unique identifier upon enrolment. We encrypted all data communication to
204 ensure data security. A team of statisticians and ultrasound examiners performed data cleaning.
205 Data cleaning included sending queries to participating centers to retrieve missing information or to
206 correct inconsistencies.

207

208 *Prediction models*

209 We assessed the diagnostic performance of subjective assessment, RMI, Benign Descriptors,
210 Simple rules, SRRisk, ADNEX with and without CA125, and the two-step strategy with and
211 without CA125. Predictions were based on information obtained at the inclusion scan and so are
212 blinded to the outcome. The results of the models were calculated after closing the study and were
213 not used to guide patient management.

214 *RMI* does not give an estimated risk but a non-negative integer (0 or higher), with higher scores
215 suggesting a higher likelihood of malignancy. It includes clinical (CA125 and menopausal status)
216 and ultrasound variables (multilocular cyst, solid areas, bilateral lesions, ascites, abdominal
217 metastases).¹⁰

218 There are *four Benign Descriptors*, which classify tumors as benign: 1) a unilocular cyst with
219 ground glass echogenicity and largest diameter <10 cm in a premenopausal woman is suggestive of
220 endometrioma (BD 1); 2) a unilocular cyst with mixed echogenicity, acoustic shadows and largest
221 diameter <10 cm in a premenopausal woman is suggestive of benign teratoma (BD 2); 3) a
222 unilocular cyst with anechoic cyst fluid, smooth internal walls and largest diameter <10 cm in a pre-
223 or post-menopausal woman is suggestive of a simple cyst or cystadenoma (BD3); 4) all other
224 unilocular cysts with smooth internal walls and largest diameter <10 cm in a pre- or post-
225 menopausal woman (BD 4) are suggestive of a benign cyst.¹⁷

226 *Simple Rules* classify tumors as benign, inconclusive, or malignant based on the presence of five
227 benign ultrasound features (unilocular cyst, smooth multilocular cyst with largest diameter <100

228 mm, acoustic shadows, solid component(s) are present but largest diameter is <7 mm, no
229 vascularization on color Doppler) and five malignant ultrasound features (irregular solid tumor,
230 irregular multilocular solid tumor with largest diameter ≥ 100 mm, at least 4 papillary projections,
231 presence of ascites, very strong vascularization on color Doppler). The classification is inconclusive
232 if none of the ten features is present, or if both benign and malignant features are present.¹⁸ We
233 added inconclusive tumors to those predicted to be malignant, resulting in a binary classifier
234 (benign or malignant).

235 *SRRisk* is a logistic regression model that calculates the risk that an adnexal tumor is malignant
236 based on type of center (oncology center vs other center) and on the ten binary ultrasound features
237 used in the Simple Rules.²⁰

238 *ADNEX* is a multinomial regression model that calculates the probability of five outcome
239 categories: benign, borderline, stage I primary invasive ovarian malignancy, stage II-IV primary
240 invasive ovarian malignancy, and metastasis from another primary tumor.²¹ One minus the
241 probability of a benign tumor equals the estimated risk of malignancy. *ADNEX* includes three
242 clinical variables (age, type of center, i.e. oncology vs non-oncology center, serum level of CA125)
243 and six ultrasound variables (maximum diameter of the lesion in mm, proportion of solid tissue
244 calculated as the maximum diameter of the largest solid component in mm divided by the maximum
245 diameter of the lesion in mm, presence of more than 10 cyst locules, number of papillary
246 projections, acoustic shadows, ascites). The variable CA125 is optional, but for multinomial
247 discrimination *ADNEX* works better with than without CA125.²²

248 The *two-step strategy* uses the Benign Descriptors as a first step. If a Benign Descriptor applies, the
249 mass is classified as benign, if not, *ADNEX* with or without CA125 is used to estimate the risk of
250 malignancy.¹⁷ Information about how we estimated the risk of malignancy when a Benign
251 Descriptor applied is found in Appendix 2.

252

253 *Statistical analysis and sample size*

254 We followed a prespecified statistical analysis plan. The statistical analyses were performed with R
255 version 4.1.2. The adequacy of the sample size is discussed in Appendix 3. We report all
256 performance measures with 95% confidence intervals (CI).

257 Despite it being strongly recommended to collect blood samples for measurement of serum CA125
258 in all patients, CA125 results were missing in some patients. Missing CA125 values were imputed.
259 We performed multiple imputation using the method of fully conditional specification, generating
260 100 imputations, leading to 100 completed datasets. To estimate the missing values of CA125, we
261 used predictive mean matching regression using the outcome and variables that are probably related
262 to either the level of CA125 itself, or to the unavailability of CA125. The multiple imputation
263 procedure is described in Appendix 4. To calculate the predictions for each model, we used the
264 formula present in the original paper.

265 Results are presented as absolute frequency (percentage) for nominal variables and as median,
266 interquartile range (IQR) and range (min-max) for continuous variables as appropriate. We report
267 the percentage of tumors to which a Benign Descriptor applied and the outcome of masses to which
268 a Benign Descriptor applied (pooled analysis).

269 We calculated center-specific area under the receiver operating characteristic curve (AUROC) to
270 estimate the ability to discriminate between benign and malignant adnexal masses for RMI and the
271 risk models (SRRisk, ADNEX, IOTA two-step strategy) and used meta-analysis to obtain the
272 overall AUROC per model. The heterogeneity between centers was assessed by calculating 95%
273 prediction intervals (PI). The meta-analysis procedure is described in Appendix 5.

274 We assessed calibration of the risk models by calculating observed over expected ratio (O:E). O:E
275 is the ratio of the observed risk of having the outcome divided by the risk estimated by the model.
276 An O:E higher than 1 indicates that the model underestimates the risk of malignancy, and an O:E
277 lower than 1 indicates that the model overestimates the risk of malignancy. The ideal value is 1.^{41,42}

278 We also constructed flexible calibration curves using loess.⁴³ We calculated center-specific O:E
279 and calibration curves and combined them using meta-analysis. For the meta-analysis of calibration

280 curves, we combined 13 centers with small sample size or low prevalence of malignancy into four
281 groups to avoid computational problems. The four groups were: Santorso, Foggia, Treviso (group
282 1); Messina, Carpi, Montebelluna (group 2); Verona, Firenze, Padova, Rome (group 3); and Bari
283 B, Asti, Bolzano (group 4). The meta-analysis procedure is described in Appendix 5.

284 Clinical utility to decide which patients to refer for specialized oncological care was estimated
285 using decision curve analysis for risk thresholds between 1% and 50%.⁴⁴ Net benefit is a measure
286 of clinical utility. To know if a model is clinically useful, we compare it to treat all and treat none
287 (in this case to refer all or to refer none to an oncology center). A model is clinically useful if it is
288 superior to both treat all and treat none. Because RMI does not provide risk estimates, for RMI we
289 computed clinical utility at the following fixed RMI scores: 200 (a threshold often used clinically
290 and recommended in several national guidelines¹¹⁻¹⁵), 250, 100 and 25. We show overall decision
291 curves calculated using meta-analysis of center-specific decision curves. The meta-analysis
292 procedure is described in Appendix 5.

293 We calculated sensitivity, specificity, positive predictive value (PPV) and negative predictive value
294 (NPV) for subjective assessment and the Simple Rules (inconclusive cases classified as malignant)
295 and for the risk models at risk of malignancy cut-offs of 1%, 3%, 5%, 10%, 15%, 20%, 25%, 30%,
296 40%, and 50%. For RMI, we report classification performance for cut-offs 25, 100, 200, and 250.
297 We calculated center-specific sensitivity and specificity and combined them using meta-analysis.⁴⁵
298 Centers with no true positive (TP) and no false positive (FP) test results at a specific threshold were
299 excluded from the meta-analysis of PPV for that threshold. Centers with no true negative (TN) and
300 no false negative (FN) at a specific threshold were excluded from the meta-analysis of NPV for that
301 threshold.

302 To estimate the multinomial performance of ADNEX and the two-step strategy, we computed the
303 Polytomous Discrimination Index (PDI) and calculated the AUROC for each pair of outcome
304 categories using the conditional risk method.⁴⁶⁻⁴⁷ To evaluate calibration, we computed O:E ratios
305 per category. Due to the small numbers in most centers, we anticipated computational problems

306 when attempting to perform meta-analysis for multinomial performance. Therefore, we used the
307 pooled dataset to estimate the multinomial performance of ADNEX and the two-step strategy.

308

309 *Subgroup analyses*

310 We calculated AUROC and O:E ratio for prespecified subgroups based on menopausal status, type
311 of center (oncology vs non-oncology), and ultrasound examiner's experience (<500 scans
312 performed, 500-5000 scans performed, >5000 scans performed; EFSUMB Level 1, Level 2, and
313 Level 3) using pooled data due to the small numbers in most centers.

314

315 **Results**

316 A total of 1567 patients were recruited from 23 Italian centers. After data cleaning and application
317 of exclusion criteria, our study population consisted of 1431 patients in 21 Italian centers (10
318 oncological and 11 non-oncological centers) (Figure 1 and Supplementary Table S1). Based on
319 histology, 995/1431 (69.5%) tumors were benign and 436/1431 (30.5%) were malignant (115/1431,
320 8.0% borderline, 263/1431, 18.4% primary invasive, 58/1431, 4.1% metastatic tumors). Tumor
321 outcome according to center is shown in Supplementary Table S1. Clinical, ultrasound and
322 histological characteristics of the study population are summarized in Table 1. Specific histological
323 diagnoses are shown in Supplementary Table S2. The median age of the patients was 52 years (IQR
324 40-62, range 18 to 88), and 745 patients (52%) were postmenopausal. The median of the maximum
325 diameter of the lesion was 69 (IQR 48-100, range 9 - 400) mm, 281 (20%) patients had bilateral
326 masses, and 120 (8%) patients had ascites. CA125 was missing in 28% (394/1431) of patients. The
327 characteristics of our study population and those of the studies in which the RMI and IOTA models
328 were developed are shown in Supplementary Table S3. Patients in our study population were more
329 frequently postmenopausal than those in the development sets of the IOTA models, the Ca125
330 values (when available) were lower, while acoustic shadowing and absent color Doppler signals
331 were more common. Multilocular cysts were less common than in the development data set of the

332 Simple Rules, and Ca125 values were higher than in the development data set of RMI. The
333 distribution of tumor outcome (i.e., prevalence of borderline, stage I primary invasive, stage II-IV
334 primary invasive, metastatic tumors) in our study population is reasonably similar to that in the
335 studies in which the IOTA models were developed.

336 The BDs applied to 328/1431 (23%) tumors, of which 325 (99%) were benign, 3/328 (1%, 95% CI
337 1-2) were borderline, and none was an invasive malignancy (Supplementary Table S4).

338 For all IOTA models (ADNEX with and without CA125, SSRisk, two step strategy with and
339 without CA125), the overall AUROC was ≥ 0.91 , while that of RMI was 0.85 (Figure 2).
340 Differences in AUROC between centers (heterogeneity) were smallest for SSRisk (Figure 2,
341 Supplementary Figure S1-S6).

342 At a risk threshold of 10% (the risk threshold recommended in an international consensus statement
343 for referring patients to an oncology center⁶), all IOTA models (SSRisk, ADNEX with and without
344 Ca125 and the two-step strategy with and without Ca125) had sensitivity >0.9 with specificity
345 ranging from 0.77 to 0.80. ADNEX and the two step strategy had the same classification
346 performance at risk threshold 10%: sensitivity 0.92 and specificity 0.80 when CA125 was included
347 as a predictor, sensitivity 0.94 and specificity 0.77 when CA125 was not included as a predictor
348 (Supplementary Table S5). At a threshold of 200, RMI had sensitivity 0.58 and specificity 0.94
349 (Supplementary Table S6). The Simple Rules were applicable in 1244/1431 (87%) tumors and had
350 sensitivity 0.90 and specificity 0.85 (inconclusive cases classified as malignant). Subjective
351 assessment had sensitivity 0.93 and specificity 0.88 (Supplementary Table S7).

352 The malignancy risk was slightly underestimated by all IOTA models (O:E ratios 1.04 -1.20) (Table
353 2). SSRisk was slightly better calibrated than other IOTA models (point estimate O:E ratio 1.04;
354 95% CI 0.97-1.12) (Table 2, Figure 3). The relation between the RMI score and the observed
355 prevalence of malignancy is shown in supplementary Figure S7. At RMI score 200, the observed
356 prevalence of malignancy was 55% (95% CI 49-61).

357 Decision curves are shown in Figure 4. Irrespective of which RMI cutoff was used, the IOTA
358 methods had higher net benefit than RMI. At risk thresholds lower than 15%, ADNEX and the two-
359 step strategy had the highest net benefit, at risk thresholds above 20% SRRisk and subjective
360 assessment had highest net benefit. RMI at cutoff 200 was less clinically useful than simply treating
361 everyone (referring everyone to an oncology center) at risk thresholds below 15%.

362 The ability of ADNEX and the two-step strategy to discriminate between different tumor types is
363 shown in Table 3 and Supplementary Table S8. The PDI ranged from 0.44 to 0.55 for all four
364 models. ADNEX and the two-step strategy had the same ability to discriminate between the five
365 tumor types. Performance was poorest for discrimination between Stage I primary invasive tumors
366 and metastases (AUROC 0.70 for all four models) and between Stage II-IV primary invasive
367 tumors and metastases (AUROC 0.59 for models without Ca125 and 0.74 for models with CA125).
368 AUROCs ranged from 0.93 to 0.98 for distinguishing benign tumors from invasive malignancies
369 (primary or metastases). Calibration of the estimated risks for the five tumor types was similar for
370 ADNEX and the two-step strategy irrespective of whether or not CA125 was included as a predictor
371 (Supplementary Table S9). All four models underestimated the risk of borderline, stage I primary
372 invasive, and secondary metastatic tumors (O:E >1) but overestimated the likelihood of a stage II-
373 IV ovarian malignancy and benign tumor (O:E <1).

374

375 *Subgroup analyses*

376 Tumor outcome and percentage of missing CA125 values in subgroups according to menopausal
377 status, type of center and ultrasound examiners' level of experience are shown in Table 4. In all
378 subgroups the AUROCs were higher for the IOTA models than for RMI (Figure 5). The AUROCs
379 of the IOTA models were ≥ 0.90 in all subgroups (0.90 – 0.95). They were slightly higher in pre-
380 than post-menopausal patients, in non-oncology than oncology centers, and for EFSUMB Level 3
381 ultrasound examiners than for EFSUMB Level 1 or 2 ultrasound examiners. The sensitivity of
382 subjective assessment was higher for level 3 and level 2 examiners than for level 1 examiners (0.96

383 vs 0.92 vs 0.86), but the corresponding specificity was lower (0.84 vs 0.87 vs 0.97) (Supplementary
384 Table S10). The sensitivity of Simple Rules (inconclusive cases classified as malignant) was higher
385 for level 3 than level 2 and 1 examiners (0.95 vs 0.88 vs 0.82) with the corresponding specificity
386 being lower (0.80 vs 0.85 vs 0.92).

387 Calibration of the IOTA models in the subgroups is shown in Supplementary Figure S8. In all
388 subgroups, the malignancy risk tended to be underestimated (point estimate for O:E ratio >1) by the
389 IOTA models, with less underestimation of malignancy risk in postmenopausal than premenopausal
390 patients and in oncology than non-oncology centers. In all subgroups, SRRisk had the least
391 underestimation. The 95% CIs for O:E ratios were very wide for EFSUMB level 1 examiners and for
392 examiners that had performed <500 scans at study start due to the small sample size for these
393 groups.

394

395 **Discussion**

396

397 *Principal findings*

398 SRRisk, ADNEX and the IOTA two-step strategy (with or without CA125) discriminated well
399 between benign and malignant adnexal masses and were superior to RMI when validated on a
400 national basis in 21 Italian centers by ultrasound examiners with different levels of ultrasound
401 experience. All IOTA methods had higher net benefit than RMI. SSRisk, ADNEX and the two step-
402 strategy (with or without CA125) had the highest net benefit at risk thresholds below 15%.

403

404 *Strengths and limitations*

405 Our study is the first prospective national multicenter study to validate IOTA models and to validate
406 them in the hands of ultrasound examiners with different levels of experience. We assessed the
407 diagnostic performance both in terms of discrimination, calibration and clinical utility. Limitations
408 are the small number of EFSUMB level 1 examiners and examiners that had performed <500 scans

409 at the start of the study, the small number of centers from the south of Italy (our aim was to include
410 centers homogenously distributed all over Italy), and that CA125 was missing in a substantial
411 proportion of patients (28%). We adressed the missing CA125 values using multiple imputation.
412 Using histology as reference standard can be seen both as a strength and as a limitation. The
413 strength is that using the same reference standard in all patients avoids differential verification
414 bias. The limitation is that our results might not be applicable to all adnexal masses, which include
415 also those managed with clinical and ultrasound follow-up.

416

417 *Comparison with other studies*

418 The results of our study agree well with those in other validation studies.^{9,17,23,31,48-52} The
419 discriminative performance of ADNEX in our study (AUROC 0.92 and 0.91 for ADNEX with and
420 without CA125, respectively) was similar to that reported in a meta-analysis including 17007
421 adnexal masses examined with ultrasound in different countries and settings in 47 studies (AUROC
422 0.93 both for ADNEX with and without CA125).⁴⁸ It was only slightly poorer than that in a large
423 international multicenter study conducted by the IOTA group (AUROC 0.94 both for ADNEX with
424 and without CA125) and in a large single center study conducted in a private center in Barcelona,
425 Spain (AUROC 0.95).^{22,52} The discriminative performance of the two-step stragey was also
426 slightly poorer than in two other large studies^{17,52} (AUROC 0.92 vs 0.95 when ADNEX with
427 CA125 was used as second step test; AUROC 0.91 vs 0.94 vs. 0.95 when ADNEX without CA125
428 was used as second step test). Whether the small differences in discriminative performance are
429 explained by differences in tumor characteristics (the studies cited included also patients managed
430 expectantly) or in ultrasound expertise is difficult to know. We found the discriminative ability and
431 the clinical utility of ADNEX, the two-step strategy and SRRisk to be superior to those of RMI,
432 which agrees with the results of other studies.^{17,22} Both in our study and in others, the IOTA
433 models underestimated the risk of malignancy, the best calibrated model being SRRisk, and the
434 models being better calibrated in postmenopausal than premenopausal patients.^{17,22} The Benign

435 Descriptors were applicable in a lower proportion of patients in our study than in those by Landolfo
436 et al and Pascual et al (23% vs 37% vs 77%).^{17,52} This is explained by the the other two studies
437 including also patients managed expectantly and by the study by Pascual et al being an extreme
438 low-risk population.^{17,52} The classification performance of the Simple Rules (inconclusive cases
439 classified as malignant) in our study (sensitivity 0.90, specificity 0.85) is reasonably similar to that
440 reported in published meta-analyses (sensitivity and specificity 0.93 and 0.81,⁵⁰ sensitivity and
441 specificity 0.94 and 0.76,⁵¹ sensitivity and specificity 0.93 and 0.80⁹) and in the international
442 multicenter study by van Calster et al (sensitivity and specificity 0.90 and 0.87²²).

443 We found the sensitivity of subjective assessment to be 0.93 and the specificity to be 0.88, which is
444 almost identical to the sensitivity and specificity of subjective assessment reported in a meta-
445 analysis (sensitivity 0.93, specificity 0.89).⁹ The classification performance of subjective
446 assessment is heavily dependent on the experience of the ultrasound examiner.⁵³ The performance
447 of risk calculation models should be less dependent on ultrasound skill as long as the ultrasound
448 examiner is familiar with the definitions of the variables in the models. Nonetheless, we found some
449 small differences in the discriminative and calibration performance of the IOTA models between
450 examiners with different levels of experience, with performance being slightly better in the group of
451 EFSUMB level 3 examiners. However, it is difficult to interpret these differences, because they
452 may be explained by a difference in tumor types between the groups.

453

454 *Implications in clinical practice*

455 The good performance of the IOTA models in our study, which includes also ultrasound examiners
456 with little ultrasound experience and both local, regional and university hospitals, supports that
457 IOTA models can be widely applied in clinical practice. Our findings also support the
458 recommendation by Landolfo et. al¹⁷ to use the two-step strategy. The two-step strategy had almost
459 the same discrimination and calibration performance and almost the same clinical utility as ADNEX
460 at risk thresholds up to 20% (and better clinical utility than ADNEX at the lowest risk thresholds) in

461 our study, but the two-step strategy is easier to use than ADNEX while still offering the advantage
462 of providing an estimate of the likelihood of four types of malignancy.

463

464 *Future perspectives*

465 Prospective studies including a very large number of ultrasound examiners with limited experience
466 are needed to confirm our results. It would be important to investigate the effect of using the IOTA
467 models in impact studies.⁵⁴ Such studies will show whether use of IOTA models in daily practice
468 improve decision making and, ultimately, patient outcome.

469

470 *Conclusions*

471 SRRisk, ADNEX and the two step strategy with or without CA125 had similar and good ability to
472 distinguish benign from malignant adnexal tumours in patients examined by either expert or non-
473 expert ultrasound operators in Italy, and they were all superior to RMI. Our results support the
474 recommendation by the IOTA-group to use the two-step strategy to characterize ovarian tumors.

475

476

477 **References**

- 478 1. Siegel RL, Giaquinto AN, Jemal A. Cancer statistics, 2024. *CA Cancer J Clin.*
479 2024;74(1):12-49. doi: 10.3322/caac.21820.
480
- 481 2. Woo YL, Kyrgiou M, Bryant A, Everett T, Dickinson HO. Centralisation of services for
482 gynaecological cancers - a Cochrane systematic review. *Gynecol Oncol.* 2012;126(2):286-
483 90. doi: 10.1016/j.ygyno.2012.04.012.
484
- 485 3. Engelen MJ, Kos HE, Willemse PH, et al. Surgery by consultant gynecologic oncologists
486 improves survival in patients with ovarian carcinoma. *Cancer.* 2006;106(3):589-98. doi:
487 10.1002/cncr.21616.
488
- 489 4. Vernooij F, Heintz AP, Witteveen PO, van der Heiden-van der Loo M, Coebergh JW, van
490 der Graaf Y. Specialized care and survival of ovarian cancer patients in The Netherlands:
491 nationwide cohort study. *J Natl Cancer Inst.* 2008;100(6):399-406. doi:
492 10.1093/jnci/djn033.

493
494
495
496
497
498
499
500
501
502
503
504
505
506
507
508
509
510
511
512
513
514
515
516
517
518
519
520
521
522
523
524
525
526
527
528
529
530
531
532
533
534
535
536
537

5. Froyman W, Landolfo C, De Cock B, et al. Risk of complications in patients with conservatively managed ovarian tumours (IOTA5): a 2-year interim analysis of a multicentre, prospective, cohort study. *Lancet Oncol.* 2019;20(3):448-458. doi: 10.1016/S1470-2045(18)30837-4.
6. Timmerman D, Planchamp F, Bourne T, et al. ESGO/ISUOG/IOTA/ESGE Consensus Statement on preoperative diagnosis of ovarian tumors. *Ultrasound Obstet Gynecol.* 2021;58(1):148-168. doi: 10.1002/uog.23635.
7. Valentin L, Hagen B, Tingulstad S, Eik-Nes S. Comparison of 'pattern recognition' and logistic regression models for discrimination between benign and malignant pelvic masses: a prospective cross validation. *Ultrasound Obstet Gynecol.* 2001;18(4):357-65. doi: 10.1046/j.0960-7692.2001.00500.x.
8. Timmerman D. The use of mathematical models to evaluate pelvic masses; can they beat an expert operator? *Best Pract Res Clin Obstet Gynaecol.* 2004;18(1):91-104. doi: 10.1016/j.bpobgyn.2003.09.009.
9. Meys EM, Kaijser J, Kruitwagen RF, et al. Subjective assessment versus ultrasound models to diagnose ovarian cancer: A systematic review and meta-analysis. *Eur J Cancer.* 2016 May;58:17-29. doi: 10.1016/j.ejca.2016.01.007.
10. Jacobs I, Oram D, Fairbanks J, Turner J, Frost C, Grudzinskas JG. A risk of malignancy index incorporating CA 125, ultrasound and menopausal status for the accurate preoperative diagnosis of ovarian cancer. *Br J Obstet Gynaecol.* 1990;97(10):922-9. doi: 10.1111/j.1471-0528.1990.tb02448.x.
11. <https://kunskapsbanken.cancercentrum.se/diagnoser/aggstockscancer-epitelial/vardprogram/diagnostik/>
12. https://www.legeforeningen.no/contentassets/75de8a48892f4cdc8db45fc7804d369c/benigne-ovarialcyster_2024.pdf
13. <https://static1.squarespace.com/static/5467abcce4b056d72594db79/t/58ed3a3ddb29d654c6cd634c/1491941952677/Cysteguideline.pdf>.
14. Royal College of Obstetricians and Gynaecologists. The Management of Ovarian Cysts in Postmenopausal Women. Green-top Guideline No. 34. July 2016.
15. Sundar S, Agarwal R, Davenport C, et al. Risk-prediction models in postmenopausal patients with symptoms of suspected ovarian cancer in the UK (ROCKeTS): a multicentre, prospective diagnostic accuracy study. *Lancet Oncol.* 2024 Oct;25(10):1371-1386. doi: 10.1016/S1470-2045(24)00406-6.

- 538 16. Ameye L, Timmerman D, Valentin L, et al. Clinically oriented three-step strategy for
539 assessment of adnexal pathology. *Ultrasound Obstet Gynecol.* 2012;40(5):582-91. doi:
540 10.1002/uog.11177.
541
- 542 17. Landolfo C, Bourne T, Froyman W, et al. Benign descriptors and ADNEX in two-step
543 strategy to estimate risk of malignancy in ovarian tumors: retrospective validation in IOTA5
544 multicenter cohort. *Ultrasound Obstet Gynecol.* 2023 Feb;61(2):231-242. doi:
545 10.1002/uog.26080.
546
- 547 18. Timmerman D, Testa AC, Bourne T, et al. Simple ultrasound-based rules for the diagnosis
548 of ovarian cancer. *Ultrasound Obstet Gynecol.* 2008;31(6):681-90. doi: 10.1002/uog.5365.
549
- 550 19. Timmerman D, Testa AC, Bourne T, et al. Logistic regression model to distinguish between
551 the benign and malignant adnexal mass before surgery: a multicenter study by the
552 International Ovarian Tumor Analysis Group. *J Clin Oncol.* 2005;23(34):8794-801. doi:
553 10.1200/JCO.2005.01.7632.
554
- 555 20. Timmerman D, Van Calster B, Testa A, et al. Predicting the risk of malignancy in adnexal
556 masses based on the Simple Rules from the International Ovarian Tumor Analysis group.
557 *Am J Obstet Gynecol.* 2016;214(4):424-437. doi: 10.1016/j.ajog.2016.01.007.
558
- 559 21. Van Calster B, Van Hoorde K, Valentin L, et al. Evaluating the risk of ovarian cancer before
560 surgery using the ADNEX model to differentiate between benign, borderline, early and
561 advanced stage invasive, and secondary metastatic tumours: prospective multicentre
562 diagnostic study. *BMJ.* 2014;349:g5920. doi: 10.1136/bmj.g5920.
563
- 564 22. Van Calster B, Valentin L, Froyman W, et al. Validation of models to diagnose ovarian
565 cancer in patients managed surgically or conservatively: multicentre cohort study. *BMJ.*
566 2020;370:m2614. doi: 10.1136/bmj.m2614.
567
- 568 23. Hiett AK, Sonek JD, Guy M, Reid TJ. Performance of IOTA Simple Rules, Simple Rules
569 risk assessment, ADNEX model and O-RADS in differentiating between benign and
570 malignant adnexal lesions in North American women. *Ultrasound Obstet Gynecol.*
571 2022;59(5):668-676. doi: 10.1002/uog.24777.
572
- 573 24. Jeong SY, Park BK, Lee YY, Kim TJ. Validation of IOTA-ADNEX Model in
574 Discriminating Characteristics of Adnexal Masses: A Comparison with Subjective
575 Assessment. *J Clin Med.* 2020;9(6):2010. doi: 10.3390/jcm9062010.
576
- 577 25. Esquivel Villabona AL, Rodríguez JN, Ayala N, et al. Two-Step Strategy for Optimizing the
578 Preoperative Classification of Adnexal Masses in a University Hospital, Using International
579 Ovarian Tumor Analysis Models: Simple Rules and Assessment of Different NEoplasias in
580 the adneXa Model. *J Ultrasound Med.* 2022;41(2):471-482. doi: 10.1002/jum.15728.
581
- 582 26. Rashmi N, Singh S, Begum J, Sable MN. Diagnostic Performance of Ultrasound-Based
583 International Ovarian Tumor Analysis Simple Rules and Assessment of Different

- 584 NEoplasias in the adneXa Model for Predicting Malignancy in Women with Ovarian
585 Tumors: A Prospective Cohort Study. *Womens Health Rep* (New Rochelle). 2023;4(1):202-
586 210. doi: 10.1089/whr.2022.0072.
- 587
- 588 27. Velayo CL, Reforma KN, Sicam RVG, Diwa MH, Sy ADR, Tantengco OAG. Diagnostic
589 Performances of Ultrasound-Based Models for Predicting Malignancy in Patients with
590 Adnexal Masses. *Healthcare* (Basel). 2022;11(1):8. doi: 10.3390/healthcare11010008.
- 591
- 592 28. Qian L, Du Q, Jiang M, Yuan F, Chen H, Feng W. Comparison of the Diagnostic
593 Performances of Ultrasound-Based Models for Predicting Malignancy in Patients With
594 Adnexal Masses. *Front Oncol*. 2021;11:673722. doi: 10.3389/fonc.2021.673722.
- 595
- 596 29. Araujo KG, Jales RM, Pereira PN, et al. Performance of the IOTA ADNEX model in
597 preoperative discrimination of adnexal masses in a gynecological oncology center.
598 *Ultrasound Obstet Gynecol*. 2017;49(6):778-783. doi: 10.1002/uog.15963.
- 599
- 600 30. Timmerman D, Van Calster B, Testa AC, et al. Ovarian cancer prediction in adnexal masses
601 using ultrasound-based logistic regression models: a temporal and external validation study
602 by the IOTA group. *Ultrasound Obstet Gynecol*. 2010;36(2):226-34. doi: 10.1002/uog.7636.
- 603
- 604 31. Timmerman D, Ameye L, Fischerova D, et al. Simple ultrasound rules to distinguish
605 between benign and malignant adnexal masses before surgery: prospective validation by
606 IOTA group. *BMJ*. 2010;341:c6839. doi: 10.1136/bmj.c6839.
- 607
- 608 32. Sayasneh A, Ferrara L, De Cock B, et al. Evaluating the risk of ovarian cancer before
609 surgery using the ADNEX model: a multicentre external validation study. *Br J Cancer*.
610 2016;115(5):542-8. doi: 10.1038/bjc.2016.227.
- 611
- 612 33. Poonyakanok V, Tanmahasamut P, Jaishuen A, et al. Preoperative Evaluation of the
613 ADNEX Model for the Prediction of the Ovarian Cancer Risk of Adnexal Masses at Siriraj
614 Hospital. *Gynecol Obstet Invest*. 2021;86(1-2):132-138. doi: 10.1159/000513517.
- 615
- 616 34. Giourga M, Pouliakis A, Vlastarakos P, et al. Evaluation of IOTA-ADNEX Model and
617 Simple Rules for Identifying Adnexal Masses by Operators with Varying Levels of
618 Expertise: A Single-Center Diagnostic Accuracy Study. *Ultrasound Int Open*.
619 2023;9(1):E11-E17. doi: 10.1055/a-2044-2855.
- 620
- 621 35. Grover SB, Patra S, Grover H, Mittal P, Khanna G. Prospective revalidation of IOTA "two-
622 step", "alternative two-step" and "three-step" strategies for characterization of adnexal
623 masses - An Indian study focussing the radiology context. *Indian J Radiol Imaging*.
624 2020;30(3):304-318. doi: 10.4103/ijri.IJRI_279_20.
- 625
- 626 36. <https://iotaplus.org/en/certified-members>.
- 627
- 628 37. European Federation of Societies for Ultrasound in Medicine and Biology. *Ultraschall Med*
629 2006;27(1): 79-95 doi: 10.1055/s-2006-933605

630
631
632
633
634
635
636
637
638
639
640
641
642
643
644
645
646
647
648
649
650
651
652
653
654
655
656
657
658
659
660
661
662
663
664
665
666
667
668
669
670
671
672
673

38. Timmerman D, Valentin L, Bourne TH, Collins WP, Verrelst H, Vergote I; International Ovarian Tumor Analysis (IOTA) Group. Terms, definitions and measurements to describe the sonographic features of adnexal tumors: a consensus opinion from the International Ovarian Tumor Analysis (IOTA) Group. *Ultrasound Obstet Gynecol*. 2000;16(5):500-5.
39. Prat J; FIGO Committee on Gynecologic Oncology. Staging classification for cancer of the ovary, fallopian tube, and peritoneum. *Int J Gynaecol Obstet*. 2014;124(1):1-5. doi: 10.1016/j.ijgo.2013.10.001.
40. Installé AJ, Van den Bosch T, De Moor B, Timmerman D. Clinical data miner: an electronic case report form system with integrated data preprocessing and machine-learning libraries supporting clinical diagnostic model research. *JMIR Med Inform*. 2014;2(2):e28. doi: 10.2196/medinform.3251.
41. Riley RD, Debray TPA, Collins GS, et al. Minimum sample size for external validation of a clinical prediction model with a binary outcome. *Stat Med*. 2021;40(19):4230-4251. doi: 10.1002/sim.9025.
42. Debray TP, Damen JA, Snell KI, et al. A guide to systematic review and meta-analysis of prediction model performance. *BMJ*. 2017;356:i6460. doi: 10.1136/bmj.i6460.
43. Van Calster B, Nieboer D, Vergouwe Y, De Cock B, Pencina MJ, Steyerberg EW. A calibration hierarchy for risk models was defined: from utopia to empirical data. *J Clin Epidemiol*. 2016;74:167-76. doi: 10.1016/j.jclinepi.2015.12.005.
44. Vickers AJ, Elkin EB. Decision curve analysis: a novel method for evaluating prediction models. *Med Decis Making*. 2006;26(6):565-74. doi: 10.1177/0272989X06295361.
45. Lin L, Chu H. Meta-analysis of Proportions Using Generalized Linear Mixed Models. *Epidemiology*. 2020;31(5):713-717. doi: 10.1097/EDE.0000000000001232.
46. Van Calster B, Vergouwe Y, Looman CW, Van Belle V, Timmerman D, Steyerberg EW. Assessing the discriminative ability of risk models for more than two outcome categories. *Eur J Epidemiol*. 2012;27(10):761-70. doi: 10.1007/s10654-012-9733-3.
47. Van Calster B, Van Belle V, Vergouwe Y, Timmerman D, Van Huffel S, Steyerberg EW. Extending the c-statistic to nominal polytomous outcomes: the Polytomous Discrimination Index. *Stat Med*. 2012;31(23):2610-26. doi: 10.1002/sim.5321.
48. Barreñada L, Ledger A, Dhiman P, et al. ADNEX risk prediction model for diagnosis of ovarian cancer: systematic review and meta-analysis of external validation studies. *BMJ Med*. 2024;3(1):e000817. doi: 10.1136/bmjmed-2023-000817.

- 674 49. Alcázar JL, Pascual MA, Graupera B, et al. External validation of IOTA simple descriptors
675 and simple rules for classifying adnexal masses. *Ultrasound Obstet Gynecol.*
676 2016;48(3):397-402. doi: 10.1002/uog.15854.
677
- 678 50. Kaijser J, Sayasneh A, Van Hoorde K, et al. Presurgical diagnosis of adnexal tumours using
679 mathematical models and scoring systems: a systematic review and meta-analysis. *Hum*
680 *Reprod Update.* 2014;20(3):449-62. doi: 10.1093/humupd/dmt059.
681
- 682 51. Westwood M, Ramaekers B, Lang S, et al. Risk scores to guide referral decisions for people
683 with suspected ovarian cancer in secondary care: a systematic review and cost-effectiveness
684 analysis. *Health Technol Assess.* 2018;22(44):1-264. doi: 10.3310/hta22440.
685
- 686 52. Pascual MA, Vancraeynest L, Timmerman S, et al. Validation of ADNEX and IOTA two-
687 step strategy and estimation of risk of complications during follow-up of adnexal masses in
688 low-risk population. *Ultrasound Obstet Gynecol.* 2024;64(3):395-404. doi:
689 10.1002/uog.27642.
690
- 691 53. Timmerman D, Schwärzler P, Collins WP, et al. Subjective assessment of adnexal masses
692 with the use of ultrasonography: an analysis of interobserver variability and experience.
693 *Ultrasound Obstet Gynecol.* 1999;13(1):11-6. doi: 10.1046/j.1469-0705.1999.13010011.x.
694
- 695 54. Moons KG, Altman DG, Vergouwe Y, Royston P. Prognosis and prognostic research:
696 application and impact of prognostic models in clinical practice. *BMJ.* 2009;338:b606. doi:
697 10.1136/bmj.b606.
698
699
700
701
702
703
704
705
706
707
708
709
710
711
712
713
714
715
716
717
718

719
720
721
722
723
724
725
726
727
728
729
730
731
732
733
734
735

Tables

736
737
738

Table 1. Clinical, ultrasound, and histological characteristics of the study population (n=1431).

	<i>Median (IQR), or n (%), Range</i>
Patient age at recruitment (years)	52 (IQR 40-62) Range: 18-88
Postmenopausal	745 (52)
Gynecological symptoms during the year preceding inclusion	601 (42)
Bilateral masses	281 (20)
Presence of solid components	730 (51)
Maximum diameter of lesion (mm)	69 (IQR 48-100) Range: 9-400
Largest diameter of largest solid component (mm)*	40 (IQR 16-68) Range: 2-250
Number of papillary projections	
0	1086 (76)
1	153 (11)
2	49 (3)
3	28 (2)
>3	115 (8)
More than 10 cyst locules	160 (11)
Acoustic shadows	327 (23)
Ascites	120 (8)
CA125 results missing	394 (28)
CA125 (U/mL, if available)	19 (IQR 10-57) Range: 1-12000
Colour score of intratumoural flow	
1: no blood flow	733 (51)
2: minimal blood flow	277 (19)
3: moderate blood flow	239 (17)

4: very strong blood flow	182 (13)
Histological diagnosis†	
Benign	995 (69)
Borderline	115 (8)
Stage I primary invasive	109 (8)
Stage II-IV primary invasive	154 (11)
Secondary metastatic	58 (4)

739
740
741
742
743
744
745
746
747
748
749
750
751
752
753

IQR: Interquartile range
*For tumors with a solid component
† specific histological diagnoses are shown in Supplementary Table S2.

Table 2. Calibration in terms of observed over expected ratio based on meta-analysis of data from

Model	O:E ratio (95% CI)
SRRisk	1.04 (0.97; 1.12)
ADNEX without CA125	1.11 (1.04; 1.20)
ADNEX with CA125	1.18 (1.07; 1.29)
Two-step without CA125	1.13 (1.05; 1.21)
Two-step with CA125	1.20 (1.09; 1.32)

754
755
756
757
758
759
760
761
762
763
764
765
766
767
768
769
770
771
772
773
774
775

O:E ratio, observed over expected ratio. Measure of calibration in the large (mean calibration) is calculated as the observed risk of having the outcome event in the entire validation dataset divided by the average risk predicted by the model. A value >1 indicates that the model is underestimating the average risk. A value <1 means that the model is over-estimating the average risk. CI; Confidence Interval, SRRisk; Simple Rules Risk Model, ADNEX; Assessment of Different NEoplasias in the adnexa.

776
777
778
779
780
781
782
783
784
785
786
787
788
789
790
791
792
793
794
795
796
797
798
799
800
801
802
803
804
805
806
807
808
809
810
811
812
813
814
815
816

Table 3. Polytomous discrimination index (PDI) of Assessment of Different NEoplasias in the adneXa (ADNEX), and of the two-step strategy (pooled analysis).

	PDI (95% CI)
ADNEX without CA125	0.49 (0.47; 0.53)
ADNEX with CA125	0.55 (0.51; 0.59)
Two-step strategy without CA125	0.49 (0.47; 0.52)
Two-step strategy with CA125	0.55 (0.51; 0.59)

817
818
819
820
821
822
823
824
825
826
827
828
829
830
831
832
833
834
835
836
837
838

Table 4. Tumor outcome and percentage of missing CA 125 values for all pre-specified subgroups

Subgroup	N	Outcome		Missing CA125
		Benign	Malignant	
Postmenopausal	745	444 (60)	301 (40)	176 (24)
Premenopausal	686	551 (80)	135 (20)	218 (32)
Oncology center	817	583 (64)	331 (36)	219 (27)
Other center	614	412 (80)	105 (20)	175 (28)
Level of experience of ultrasound examiners				
<500 scans	123	102 (83)	21 (17)	41 (33)
500-5000 scans	650	476 (73)	174 (27)	108 (17)
>5000 scans	658	417 (63)	241 (37)	245 (37)
EFSUMB level of ultrasound examiners				
Level 1	118	96 (81)	22 (19)	38 (32)
Level2	884	605 (68)	279 (32)	248 (28)
Level 3	429	294 (69)	135 (31)	108 (25)

839 Results are shown as n (%)
840 EFSUMB; European Federation of Societies for Ultrasound in Medicine and Biology

841
842
843
844
845
846
847
848
849
850
851

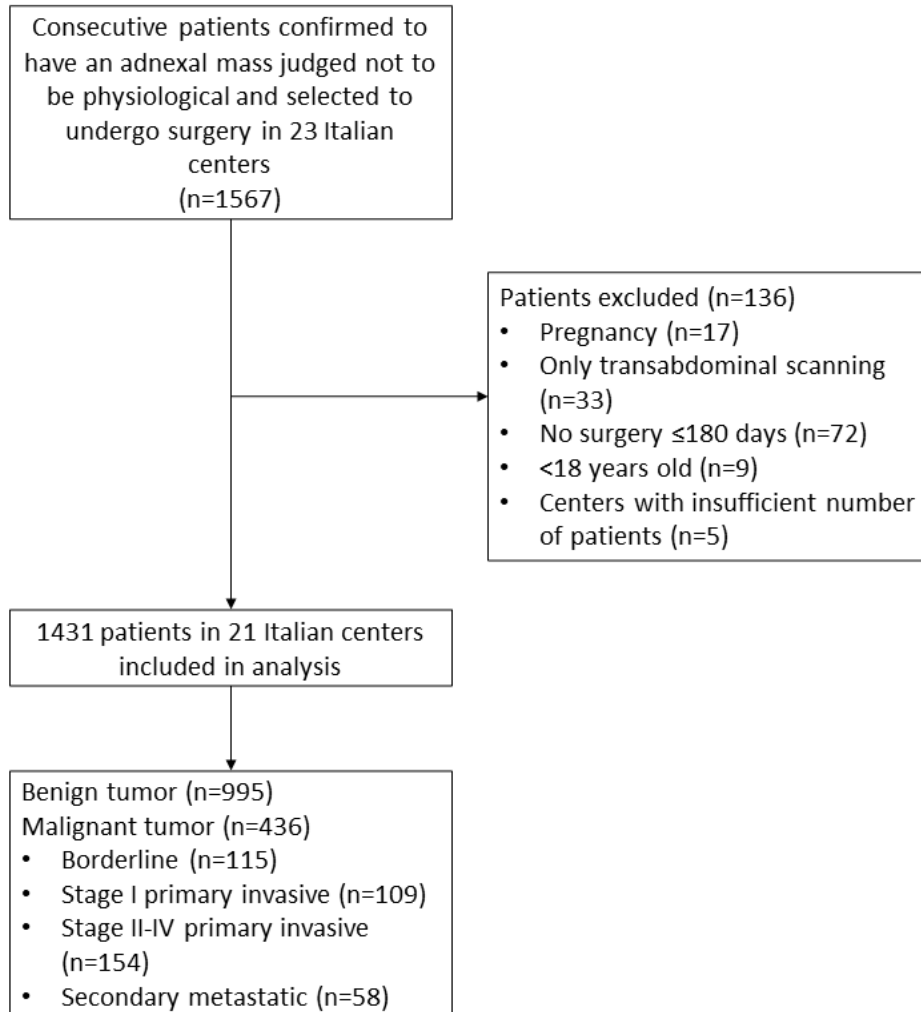
852

853

854

855

856 **Figure 1.** Flowchart of the patients included for the analysis.

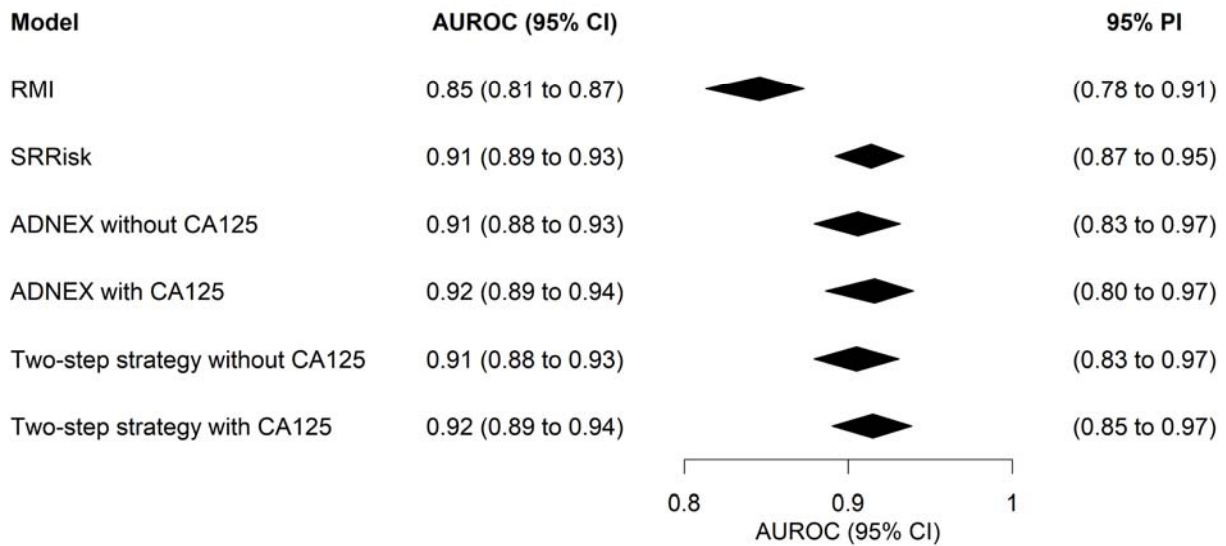


857

858

859

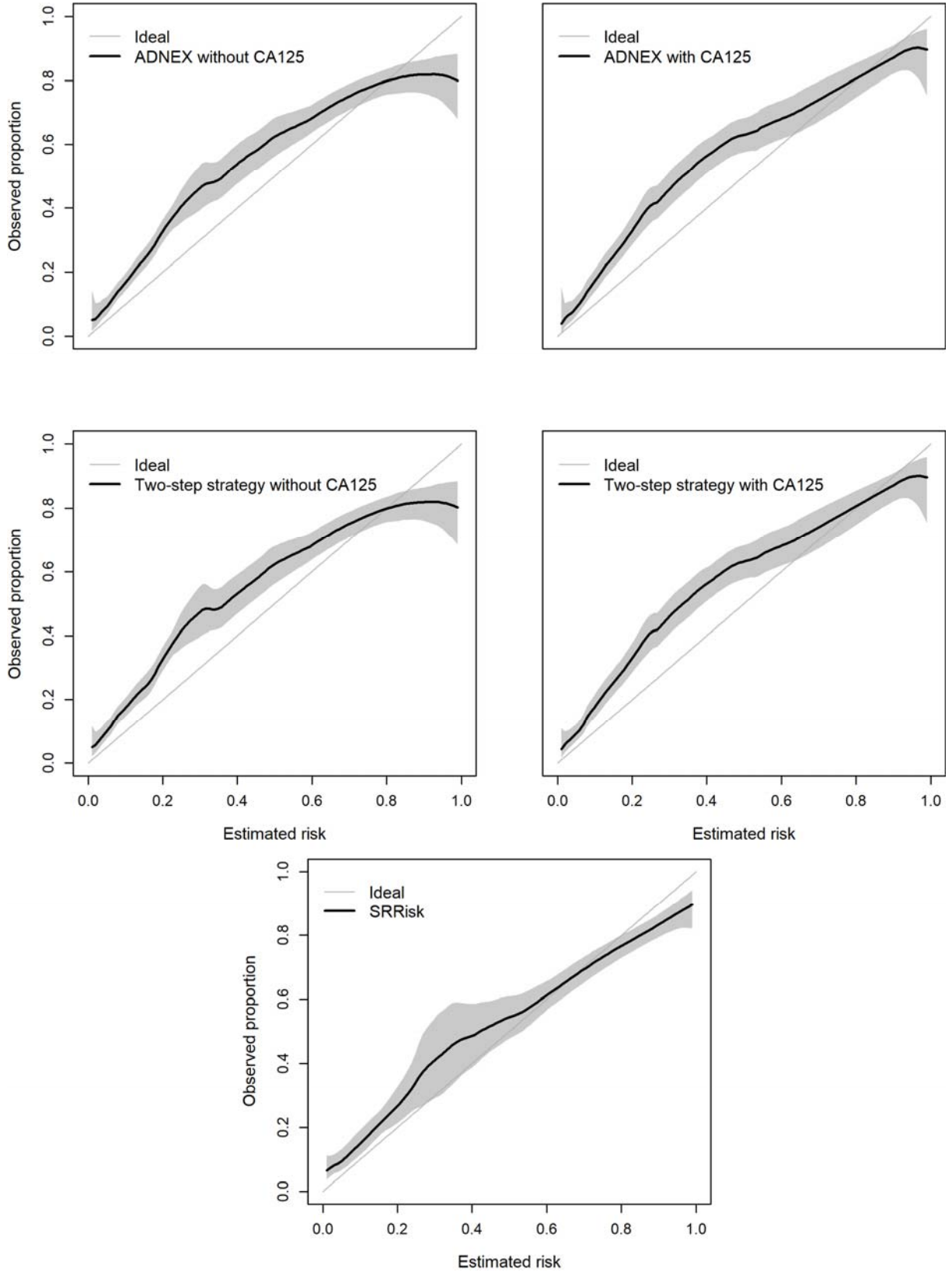
860 **Figure 2.** Summary forest plot of area under the receiver operating characteristic curve (AUROC)
861 based on meta-analysis of data from 21 centers. CI; Confidence Interval, PI; Prediction Interval.



862
863
864

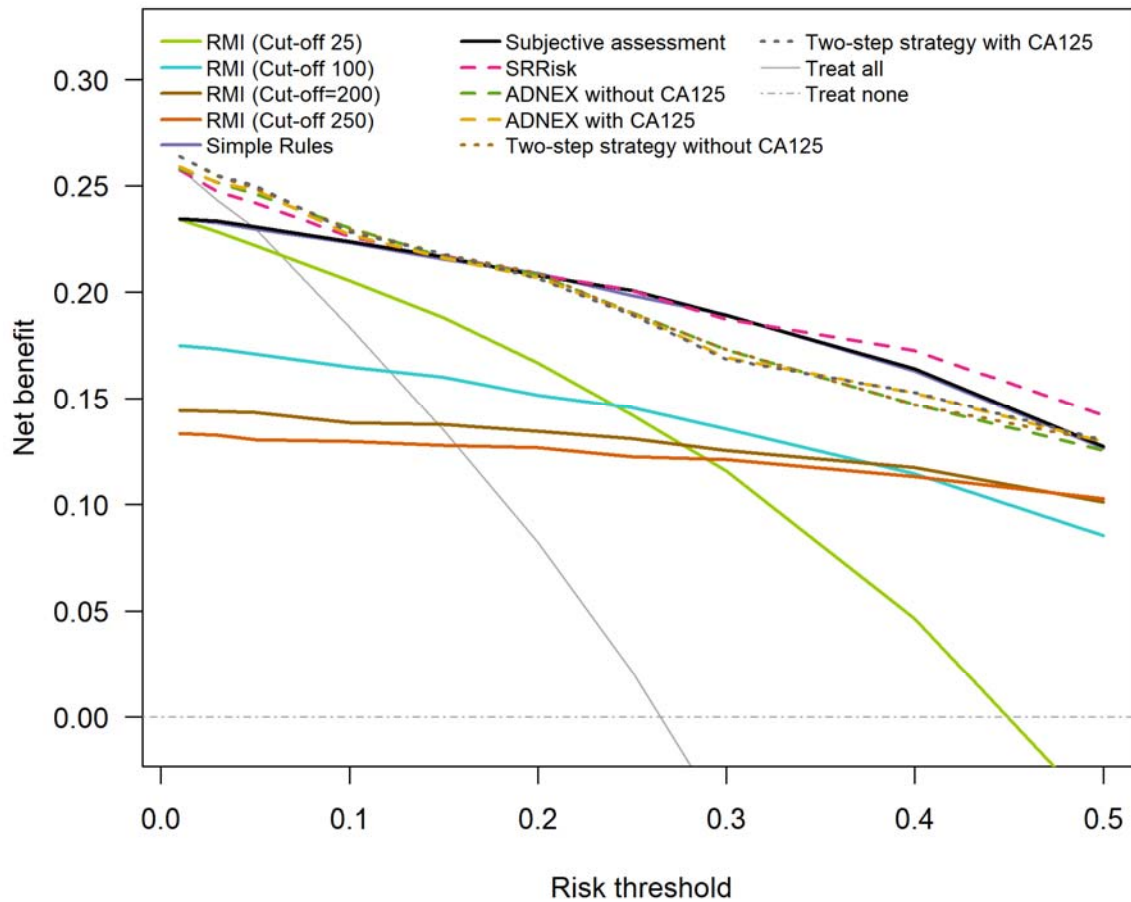
865 **Figure 3.** Flexible calibration curves using loess based on meta-analysis. Due to computational
866 problems, we divided 13 centers with low sample size or low prevalence of malignancy into four
867 groups: Santorso, Foggia, Treviso (group 1); Messina, Carpi, Montebelluna (group 2); Verona,
868 Firenze, Padova, Roma (group 3); Bari B, Asti, Bolzano (group 4). The curves were obtained with
869 meta-analysis of center-specific curves from eight centers and from the four groups. SRRisk;
870 Simple Rules Risk Model, ADNEX; Assessment of Different NEoplasias in the adnexa.

It is made available under a [CC-BY 4.0 International license](https://creativecommons.org/licenses/by/4.0/).



871
872
873

874 **Figure 4.** Decision curves for risk models, Risk of Malignancy Index (RMI), Simple Rules and
875 subjective assessment based on meta-analysis of data from 21 centers. The curves show net benefit
876 at several thresholds between 1% and 50%. A model is clinically useful if it is superior to both treat
877 all and treat none. At risk thresholds below 15%, using RMI at cutoff 200 is worse than treating
878 everyone (i.e., worse than referring all women with an adnexal mass to a gynecological oncology
879 center).



880
881
882

883 **Figure 5.** Forest plot of area under the receiver operating characteristic curve (AUROC) for
 884 prespecified subgroups (pooled data). CI; Confidence interval, PI; Prediction interval, EFSUMB;
 885 European Federation of Societies for Ultrasound in Medicine and Biology.
 886

