



LUND UNIVERSITY

Strategies to diagnose ovarian cancer: new evidence from phase 3 of the multicentre international IOTA study.

Testa, A; Kaijser, J; Wynants, L; Fischerova, D; Van Holsbeke, C; Franchi, D; Savelli, L; Epstein, Elisabeth; Czekierdowski, A; Guerriero, S; Fruscio, R; Leone, F P G; Vergote, I; Bourne, T; Valentin, Lil; Van Calster, B; Timmerman, D

Published in:
British Journal of Cancer

DOI:
[10.1038/bjc.2014.333](https://doi.org/10.1038/bjc.2014.333)

2014

[Link to publication](#)

Citation for published version (APA):

Testa, A., Kaijser, J., Wynants, L., Fischerova, D., Van Holsbeke, C., Franchi, D., Savelli, L., Epstein, E., Czekierdowski, A., Guerriero, S., Fruscio, R., Leone, F. P. G., Vergote, I., Bourne, T., Valentin, L., Van Calster, B., & Timmerman, D. (2014). Strategies to diagnose ovarian cancer: new evidence from phase 3 of the multicentre international IOTA study. *British Journal of Cancer*, 111(4), 680-688.
<https://doi.org/10.1038/bjc.2014.333>

Total number of authors:
17

General rights

Unless other specific re-use rights are stated the following general rights apply:
Copyright and moral rights for the publications made accessible in the public portal are retained by the authors and/or other copyright owners and it is a condition of accessing publications that users recognise and abide by the legal requirements associated with these rights.

- Users may download and print one copy of any publication from the public portal for the purpose of private study or research.
- You may not further distribute the material or use it for any profit-making activity or commercial gain
- You may freely distribute the URL identifying the publication in the public portal

Read more about Creative commons licenses: <https://creativecommons.org/licenses/>

Take down policy

If you believe that this document breaches copyright please contact us providing details, and we will remove access to the work immediately and investigate your claim.

LUND UNIVERSITY

PO Box 117
221 00 Lund
+46 46-222 00 00

Strategies to diagnose ovarian cancer: new evidence from phase 3 of the multicentre international IOTA study.

Running head title: Diagnosing ovarian cancer: evidence from IOTA3

Antonia Testa^{1,†}, Jeroen Kaijser^{2,3,†}, Laure Wynants^{4,5}, Daniela Fischerova⁶, Caroline Van Holsbeke⁷, Dorella Franchi⁸, Luca Savelli⁹, Elisabeth Epstein¹⁰, Artur Czekierdowski¹¹, Stefano Guerriero¹², Robert Fruscio¹³, Francesco Paolo Giuseppe Leone¹⁴, Ignace Vergote³, Tom Bourne^{2,3,15}, Lil Valentin¹⁷, Ben Van Calster², Dirk Timmerman^{2,3}

[†] “The first two authors should be regarded as joint First Authors”

¹ Department of Gynaecologic Oncology, Catholic University of the Sacred Heart, Largo Francesco Vito 8, Rome, 00165, Italy; ² KU Leuven Department of Development and Regeneration, Herestraat 49 box 7003, 3000 Leuven, Belgium; ³ Department of Obstetrics and Gynaecology and Leuven Cancer Institute, University Hospitals Leuven, Herestraat 49 box 7003, 3000 Leuven, Belgium; ⁴ KU Leuven Department of Electrical Engineering (ESAT-STADIUS), Kasteelpark Arenberg 10, 3001 Leuven, Belgium; ⁵ iMinds Future Health Department, KU Leuven, Kasteelpark Arenberg 10, 3001 Leuven, Belgium; ⁶ Gynaecological Oncology Center, Department of Obstetrics and Gynaecology, Charles University, Apolinarska 18, 120 00 Prague, Czech Republic; ⁷ Department of Obstetrics and Gynaecology, Ziekenhuis Oost Limburg, Schiepse Bos 6, 3600 Genk, Belgium; ⁸ Preventive Gynaecology Unit, Division of Gynaecology, European Institute of Oncology, via Ripamonti 435, Milan 20141, Italy; ⁹ Gynaecology and Reproductive Medicine Unit, S. Orsola-Malpighi Hospital, University of Bologna, Via Albertoni 15, Bologna 40138, Italy; ¹⁰ Department of Obstetrics and Gynaecology, Karolinska University Hospital, Solna, SE-17176 Stockholm, Sweden; ¹¹ Ist Department of Gynecologic Oncology and Gynecology, Medical University of Lublin, ul. Staszica 16, 20-081 Lublin, Poland; ¹² Department of Obstetrics and Gynaecology, Azienda Ospedaliero Universitaria di Cagliari, Strada Statale 554 Monserrato, Cagliari 09045, Italy; ¹³ Division of Obstetrics and Gynaecology, San Gerardo Hospital, University of Milan-Bicocca, Via Pergolesi, 33, Monza 20052, Italy; ¹⁴ Department of Obstetrics and Gynecology, Clinical Sciences Institute L. Sacco, University of Milan, Via G.B. Grassi 74, 20157 Milan, Italy; ¹⁵ Queen Charlotte’s and Chelsea Hospital, Imperial College, Du Cane Road, London W12 0HS, United Kingdom; ¹⁶ Department of Obstetrics and Gynaecology, Skåne University Hospital Malmö, Lund University, Södra Förstadsgatan, 20502 Malmö, Sweden;

Corresponding author: Professor Dirk Timmerman, MD PhD. Department of Obstetrics and Gynaecology and Leuven Cancer Institute. University Hospital KU Leuven. Herestraat 49 B-3000 Leuven, Belgium. T 0032 16 344201. E dirk.timmerman@uzleuven.be

Abstract

Background: To compare different ultrasound-based International Ovarian Tumour Analysis (IOTA) strategies and Risk of Malignancy Index (RMI) for ovarian cancer diagnosis using a meta-analysis approach of centre-specific data from IOTA 3.

Methods: This prospective multicentre diagnostic accuracy study included 2403 patients with 1423 benign and 980 malignant adnexal masses from 2009 until 2012. All patients underwent standardised transvaginal ultrasonography. Test performance of RMI, subjective assessment of ultrasound findings (SA), two IOTA risk models (LR1, LR2), and strategies involving combinations of IOTA Simple Rules (SR), Simple Descriptors (SD) and LR2 with and without SA was estimated using a meta-analysis approach. Reference standard was histology after surgery.

Results: The areas under the receiver operator characteristic curves of LR1, LR2, SA and RMI were 0.930 (0.917-0.942), 0.918 (0.905-0.930), 0.914 (0.886-0.936) and 0.875 (0.853-0.894). Diagnostic one and two-step strategies using LR1, LR2, SR, and SD achieved summary estimates for sensitivity 90-96%, specificity 74%-79% and diagnostic odds ratio (DOR) 32.8.-50.5. Adding SA when IOTA methods yielded equivocal results improved performance (DOR 57.6-75.7). RMI had sensitivity 67%, specificity 91% and DOR 17.5.

Conclusion: This study shows all IOTA strategies had excellent diagnostic performance in comparison to RMI. The IOTA strategy chosen may be determined by clinical preference.

Introduction

Providing care within highly specialised multi-disciplinary services has a clear survival-benefit for patients with ovarian cancer (Woo *et al*, 2012). Although such centralised care is recommended in many developed countries, a large proportion of ovarian cancer patients remain treated by general surgeons and physicians (Verleye *et al*, 2010). Several factors probably contribute to this failure to refer for specialist care, but the lack of effective preoperative strategies to evaluate ovarian tumours is certainly one of the most important (Miller and Ueland, 2012). Reports from the International Ovarian Tumour Analysis (IOTA) multicentre studies phase 1, 1b, 2 and 4 (Sayasneh *et al*, 2013a; Sayasneh *et al*, 2013b; Timmerman *et al*, 2010a; Van Holsbeke *et al*, 2012; Kaijser *et al*, 2013a; Timmerman *et al*, 2007; Timmerman *et al*, 2005; Timmerman *et al*, 2008; Ameye *et al*, 2012; Timmerman *et al*, 2010b; Van Holsbeke *et al*, 2009) have demonstrated that IOTA ultrasound-based approaches to characterise adnexal masses in the hands of physicians and sonographers with varying levels of experience outperform other established strategies such as use of individual biomarkers (serum CA-125), the Risk of Malignancy Index (RMI) (Jacobs *et al*, 1990), or Risk of Ovarian Malignancy Algorithm (ROMA) (Moore *et al*, 2009), for the classification of ovarian pathology. Nevertheless, there is a paucity of comprehensive prospective studies comparing different diagnostic strategies for ovarian cancer diagnosis on the same study population. Such studies are of pivotal importance for assessing diagnostic test accuracy. In the most recently published meta-analysis only a small number of the studies included validated different diagnostic tests for ovarian cancer on the same dataset (Kaijser *et al*, 2013b).

The primary aim of this study, the IOTA phase 3 study, was to compare the test performance of various IOTA diagnostic strategies and RMI on prospectively collected data from a large number of patients and centres.

Materials and methods

Study design. This was a multicentre cross-sectional diagnostic accuracy study with prospective data collection. Patients were recruited between October 2009 and May 2012, in 18 centres in 6 countries (Sweden, Belgium, Italy, Poland, Spain and Czech Republic). These centres were either oncology referral centres (i.e. tertiary centres for the treatment of women with gynaecological malignancy) or general hospitals and units with a special interest in gynaecological ultrasound. The centres and type of centres in IOTA 3 are listed in a supplementary appendix. All centres except three (*SSW*, *BSP*, *FIT*) had participated in at least one of the previous IOTA studies (1, 1b, or 2). Ethics approval was obtained by the Ethics Committee of the University Hospitals Leuven as main investigating centre (B32220095331/S51375) as well of the local committees of all contributing centres to IOTA3.

Inclusion criteria.

Patients were eligible if they presented with at least one adnexal mass (ovarian, para-ovarian, or tubal), underwent transvaginal ultrasound examination by a principal investigator at one of the participating centres and were then selected for surgical intervention by the managing clinician. Patients were examined following the research protocol if they gave informed consent. If more than one adnexal mass was detected, the mass with the most complex ultrasound morphology was denoted by the ultrasound examiner as the dominant mass, i.e. the one to be used for statistical analysis. If both masses had similar morphology, the largest one or the one most easily accessible by ultrasound was denoted dominant.

Exclusion criteria. Exclusion criteria were surgical removal of the mass more than 120 days after the ultrasound examination, pregnancy at scan, and data inconsistencies that persisted after final manual data checks.

Data collection. A dedicated, secure electronic data-collection system was developed for the study (IOTA 3 Study Screen; astraia Software, Munich, Germany). Patients automatically received a unique identifier. Data security was ensured by encrypting all data communication. Data integrity and completeness were ensured by client-side checks in the system supplied by astraia and final data cleaning by a group of biostatisticians and expert ultrasound examiners in Leuven, Belgium.

Ultrasound examination. All included patients underwent standardised transvaginal ultrasonography by examiners experienced in gynaecologic ultrasound (level III) (EFSUMB, 2006). High end ultrasound systems, the same or similar to those in IOTA phase 1 and 2, were used. Grey scale and colour Doppler ultrasound imaging was used to obtain information on more than 40 morphological and blood-flow variables to characterise each adnexal mass. Details on the ultrasound examination technique and the IOTA terms and definitions used to describe adnexal pathology have been published elsewhere (Timmerman *et al*, 2000). After completing the ultrasound examination, the ultrasound examiner classified each mass as benign or malignant on the basis of his/her subjective assessment of grey scale and colour or power Doppler ultrasound findings. Each mass was classified as certainly benign, probably benign, uncertain but most probably benign, uncertain but most probably malignant, probably malignant or certainly malignant. The ultrasound information was recorded prospectively in the electronic data-collection system, was locked at the time of the examination and could not be changed thereafter. Predictions of all diagnostic strategies under consideration (except subjective assessment) were obtained centrally after the conclusion of the study, and had no

role in the decision-making process. Decision-making regarding surgery for adnexal tumours was based on clinical information (such as symptoms, age, operative risk, coexisting disease, etc.) and on the clinical ultrasound report. The clinical ultrasound report was written on the basis of the results of subjective assessment.

Serum Tumor Marker. Centres were encouraged to measure the level of serum CA-125 from all patients, but the availability of this biochemical end point was not a requirement for recruitment into the study.

Diagnostic strategies. The methods and strategies prospectively compared on the IOTA 3 data set are subjective assessment, two IOTA logistic regression models, i.e. LR1 and LR2 (Timmerman *et al*, 2005), the IOTA Simple rules (Timmerman *et al*, 2008), the IOTA Simple Descriptors (Ameys *et al*, 2012) and various combinations of these, and the RMI. The IOTA methods are briefly described in **Table 1**. Details can be found in the literature (Timmerman *et al*, 2005; Timmerman *et al*, 2008; Ameys *et al*, 2012).

We evaluated five one-stage strategies, five two-stage strategies, and two three-stage strategies.

The one-stage strategies are: the use in all patients of either LR1, LR2, subjective assessment, RMI or Simple Rules (classifying all tumours where the Simple Rules yield an inconclusive result as malignant)

The two-step strategies are: Simple Rules as a first stage test and subjective assessment for tumours in which Simple Rules yield an inconclusive result; LR2 as a first stage test and subjective assessment for tumours in which LR2 yields a predicted risk of malignancy of $\geq 5\%$ but $< 25\%$ (risk of malignancy of $\geq 5\%$ but $< 25\%$ arbitrarily being taken to represent an equivocal result); Simple Descriptors as a first stage test, Simple Rules for tumours

unclassifiable by the Simple Descriptors and tumours unclassifiable by the Simple Rules classified as malignant; Simple Descriptors as a first stage test and LR2 for those tumours where the Simple Descriptors are not applicable; Simple Descriptors as a first stage test and subjective assessment for those tumours in which the Simple Descriptors are not applicable

The three-step strategies are: Simple Descriptors as a first stage test, Simple Rules for tumours in which the Simple Descriptors are not applicable, and subjective assessment for masses in which Simple Rules are inconclusive; Simple Descriptors as a first stage test, LR2 for tumours in which the Simple Descriptors are not applicable, and subjective assessment for masses in which LR2 yields a predicted risk of $\geq 5\%$ but $< 25\%$.

Reference standard. The reference standard was the histologic classification of the excised mass as malignant or benign. Histological examination was carried out at the local centre. Central pathology review was not performed because in previous IOTA studies no significant differences in reported outcomes were observed between local and central pathology reports (Timmerman *et al*, 2005). Malignant tumours were classified according to the criteria recommended by the International Federation of Gynaecology and Obstetrics (Heintz *et al*, 2003). Borderline ovarian tumours were classified as malignant. The pathologist was blinded to the prediction outcomes of the index tests being compared.

Statistical analysis

We evaluated all strategies in terms of their ability to discriminate between benign and malignant masses. For the logistic regression models LR1 and LR2 (for details see **Table 1**) and for RMI the area under the receiver-operating characteristic curve (AUC) was computed. Using the six levels of diagnostic confidence, an AUC could also be constructed for subjective assessment. For all strategies we calculated sensitivity, specificity, positive and

negative likelihood ratio (LR+ and LR-) and diagnostic odds ratio (DOR) (Deeks, 2001). To do this for LR1, LR2 and RMI, we used the cutoffs suggested in previous work (i.e. risk of malignancy $\geq 10\%$ indicating malignancy when using LR1 and LR2, and RMI > 200 indicating malignancy). To recognise that performance may differ across centres results were computed using meta-analysis techniques (Macaskill *et al*, 2010; Riley *et al*, 2008; Van Klaveren *et al*, 2014). To obtain the average AUC and DOR estimates, random effects meta-analysis was performed, using the logit of the AUC or the log of the DOR as the outcome variable. Sensitivity and specificity were modeled simultaneously using random centre effects. LR+ and LR- were computed based on the estimated average sensitivity and specificity levels. Forest plots for LR2, Simple Rules and RMI were used to present centre-specific and combined results. Subgroup analyses for RMI and the most extensively validated IOTA methods (i.e. LR2 and Simple Rules) (Sayasneh *et al*, 2013a; Sayasneh *et al*, 2013b; Timmerman *et al*, 2010a; Nunes *et al*, 2012; Alcazar *et al*, 2013; Hartman *et al*, 2012; Nunes *et al*, 2013) were performed for pre- and postmenopausal women.

For LR2 we also assessed calibration, i.e. we tested the extent to which the estimated risks of malignancy corresponded to the observed prevalence of malignancy. This was carried out by constructing parametric (logistic) calibration curves per centre (Cox, 1958; Steyerberg, 2009; Bouwmeester *et al*, 2013). RMI and Simple Rules do not provide risk estimates but comparable centre-specific curves were obtained for RMI and Simple Rules in the following manner. For RMI, analogous logistic curves were constructed to link RMI values (based on $\log(\text{RMI}+1)$) to observed risks. For Simple Rules, the proportion of malignant masses was calculated for each classification level (benign, inconclusive, malignant).

CA-125 is not a mandatory variable in the IOTA studies and by consequence information on CA125 was missing in 40% of the patients. It is most likely that missing values mainly arose

when investigators did not consider CA-125 measurement necessary given the clinical situation and the ultrasound appearance of the mass. We used multiple imputation to handle the missing values (Sterne *et al*, 2009). We used all patients from phases 1, 1b, 2 and 3 for the imputation analysis. The method is described in more detail in a **supplementary appendix** and elsewhere (Van Calster *et al*, 2011).

Calculations were performed using SAS 9.3 (SAS Institute, Cary, USA). Forest plots were created in R (www.r-project.org) using the rmeta package.

When writing this paper we used the Standards for Reporting of Diagnostic Accuracy (STARD) guidelines (Bossuyt *et al*, 2003).

Results

In total 2541 women with adnexal masses were enrolled in this study. 138 women were excluded from the final dataset. Reasons for exclusion were: an interval of > 120 days between ultrasonography and surgery (n=66), pregnancy (n=31), data errors which could not be solved by contacting the respective principal investigators (n=28), and incomplete final histology (n=13). The final dataset included 2403 patients with 1423 (59%) benign and 980 (41%) malignant adnexal masses. There were 1049 postmenopausal patients (44%) and 1354 (56%) premenopausal patients. The prevalence of malignancy was 28% (378/1354) in premenopausal patients and 57% (602/1049) in postmenopausal patients.

The types of benign and malignant tumours based on histology and FIGO staging in the final dataset are presented in **Table 2**. The most common benign diagnoses were endometrioma, serous cystadenoma, and teratoma. There were 633/2403 (26.3%) primary invasive ovarian cancers, 153/2403 (6.4%) borderline ovarian tumours, 126/2403 (5.2%) metastatic cancers in the ovaries, and 68/2403 (2.8%) rare primary invasive ovarian malignancies (e.g. granulosa

cell tumour, Sertoli-Leydig cell tumour or dysgerminoma). Descriptive statistics for the variables included in LR1, LR2 and the Simple Rules in benign and malignant adnexal masses are shown in **Table 3**.

The test performance of the IOTA diagnostic strategies, subjective assessment and RMI when using a meta-analysis approach on centre-specific data is presented in **Table 4**. The logistic regression models LR1 (AUC 0.930; 0.917-0.942) and LR2 (AUC 0.918; 0.905-0.930) had diagnostic performance similar to expert subjective assessment (AUC 0.914; 0.886-0.936). The AUC of RMI was 0.875 (0.853-0.894). The IOTA risk models LR1 and LR2, and strategies using various combinations of Simple Rules, Simple Descriptors and LR2 achieved sensitivity 90-96% and specificity 74%-79% (**Table 4**). When expert subjective assessment was used in case Simple Rules or Simple Descriptors or both yielded an inconclusive result, or in the event that LR2 gave a risk $\geq 5\%$ but $<25\%$, specificity increased from 74%-79% to 85%-89% with slightly reduced sensitivity in most instances (**Table 4**). The sensitivity of RMI was 67% and the specificity 91%. LR2 yielded a risk of $\geq 5\%$ but $<25\%$ in 419 of the 2403 patients (17%), and 108 (26%) of these women had a malignant adnexal mass.

The IOTA Simple Rules were applicable in 1846 patients (76.8%) and could be applied slightly more frequently in premenopausal 1055/1354 (77.9%) than postmenopausal women 791/1049 (75.4%). In total 1090 tumours were classified as benign by the Simple Rules, and this was correct in 1044 cases (95.8%), 756 tumours were classified as malignant by the Simple Rules, and this was correct in 674 cases (89.2%). When the Simple Rules were applicable they achieved a sensitivity of 94% (674/720) and a specificity of 93% (1044/1126). The malignancy rate among tumours where the Simple Rules yielded an inconclusive result was 46.7% (260/557). A strategy that used Simple Rules as a first stage test and classified all inconclusive cases as malignant yielded a sensitivity of 95% (95%CI 93-97%) and a

specificity of 74% (95% CI 68-80%) (**Table 4**). Using subjective assessment by an expert examiner when Simple Rules yielded an inconclusive result lowered the sensitivity from 95% to 92% (89-94%) but increased the specificity from 74% to 89% (85-92%) (**Table 4**).

The IOTA Simple Descriptors could be applied in 1014 (42.2%) masses. The Simple Descriptors classified 549 tumours (23% of all tumours in the study) as benign of which two (0.4%) turned out to be malignant. The two misclassified malignancies were stage I borderline tumours. The Simple Descriptors classified 465 tumours (19% of all tumours in the study) as malignant, and 430 (92.5%) proved to be so. The 35 benign tumours misclassified as malignant by the Simple Descriptors consisted of 11 serous cystadenomas, ten fibromas, seven mucinous cystadenomas, four rare benign tumours, two teratomas, and one functional cyst. A total of 1389 (58%) tumours could not be categorized with the Simple Descriptors. The Simple Rules could be applied in 66% (912/1389) of the tumours unclassifiable by the descriptors. The combination of Simple Descriptors with Simple Rules characterized 80% (1926/2403) of all masses as benign or malignant. When a three-step strategy was applied (Simple Descriptors as first stage test, Simple Rules in tumours unclassifiable by the descriptors, and subjective assessment for masses in which the Simple Rules were inconclusive), sensitivity and specificity were 93% (95% CI 90-95%) and 88% (95% CI 84-91%), respectively (**Table 4**).

Table 5 shows the test performance of LR2, Simple Rules and RMI in pre- and postmenopausal patients when using a meta-analysis approach of centre-specific data. In both pre- and postmenopausal patients the IOTA strategies had higher sensitivity and lower specificity than RMI. The use of subjective assessment for masses not classifiable by the Simple Rules appeared to resolve the differences in specificity.

The sensitivity and specificity of LR2, Simple rules and RMI for histological subtypes of malignant disease and the absolute number of false-negative results for histological subtypes of malignancy are presented in **Supplementary Tables S1-S4**. The sensitivity with regard to borderline tumours, FIGO stage I invasive cancer, and metastatic disease was much higher for the main IOTA approaches than for RMI, and the AUCs for LR2 were larger than those for RMI for these subtypes of malignancy.

Figure 1 and **Supplementary Figure S1** illustrate the variation in number of included masses, prevalence of malignancy, and inter-centre differences in test performance (sensitivity and specificity) for LR2, Simple Rules combined with subjective expert assessment, Simple Rules and classifying inconclusive tumours as malignant, and RMI. The malignancy rate varied between 0% and 69%, whereas the number of enrolled cases per centre ranged from six to 443. For LR2 and Simple Rules differences between centres in sensitivity were smaller than differences in specificity whereas the inverse held true for RMI. Both IOTA methods had a higher sensitivity for cancer than RMI, irrespective of the prevalence of malignancy. Discrimination (AUCs) for LR2 was consistent in both oncology and non-oncology centres with a few exceptions for centres that enrolled a very small number of cases (**Supplementary Figure S2**). Discrimination for RMI showed some variation between the centres (**Supplementary Figure S3**). The summary estimates of test performance of LR2, Simple Rules, and RMI were similar irrespective of whether it was estimated using pooled data or meta-analysis. However, pooling underestimates uncertainty, while uncertainty is appropriately addressed by adopting meta-analysis techniques.

Figure 2 shows the results of calibration for LR2, RMI, and Simple Rules for the nine centres that contributed the largest number of patients. Calibration results differed between centres. This means that for a specific prediction from a diagnostic test (LR2 risk, RMI value or

Simple Rules category) the observed prevalence of malignancy varies between centres. For LR2, the risk of malignancy was underestimated in seven of the nine centres (calibration curves above the diagonal), slightly overestimated in one centre, and perfectly calibrated in one centre. The prevalence of malignancy in women with a RMI score of 200 varied between centres from 30 to 70%, and RMI values of 200 or more were associated with high malignancy rates.

Discussion

This comparison of IOTA risk prediction models and diagnostic strategies in different clinical environments using a meta-analysis approach showed excellent test performance for all IOTA methods to characterise adnexal masses before surgery. All IOTA strategies manifested better discrimination than RMI. Additionally, we have demonstrated inter-centre differences in test performance and calibration for LR2, Simple Rules and RMI. Our use of meta-analysis techniques to summarize data did not meaningfully change the summary measures of performance from those obtained with a standard pooled analysis but gave wider confidence intervals properly reflecting the uncertainty caused by differences between centres.

The strengths of this report include the use of a rigorous prospective ultrasound protocol with agreed terms, measurement techniques, and definitions; the use of advanced statistical methods to synthesise multicentre data and report summary estimates of test accuracy and calibration, thereby minimising the risk that results are overly influenced by a single centre recruiting many more patients than others; and the large number of patients, the many participating centres, and the different types of participating centres making our results highly likely to be generalisable. A limitation of our study is that the Simple rules, Simple Descriptors and the two-step and three-step strategies were not directly applied when scanning the patients. Instead, more than 40 clinical and ultrasound variables were

prospectively collected from each patient and later incorporated in the Simple Rules or Descriptors, or synthesised to become descriptors in the Simple Descriptors or features in the Simple Rules. Whilst this might not have influenced the performance of subjective assessment or of LR1 or LR2, it could have affected that of the other tests and of the two-step and three-step strategies. A second limitation is that information on CA125 was missing in 40% of cases. We solved this by using multiple imputation (Sterne *et al*, 2009; Van Calster *et al*, 2011). Two sensitivity analyses using only complete cases for CA125 confirmed the difference in test performance in favour of the IOTA methods (LR2 and Simple Rules) (**Supplementary Tables S5-S6**). However, these approaches are biased because CA125 is more often missing in tumours that are easy to diagnose and more likely to be benign (**Supplementary Table S7**). This explains why model performance was slightly poorer for all methods when they were tested only in cases with available CA125 results. For this reason, multiple imputation is generally considered a more appropriate method to deal with missingness than to analyse only data with complete information (Sterne *et al*, 2009). A third limitation is that most of the patients in the IOTA phase 3 study were scanned by the same experienced examiners as in the centres where the IOTA methods were developed, or by examiners that had already adopted the IOTA examination technique and terminology. This may explain why the results of IOTA phase 3 confirm those of previous IOTA studies that showed excellent test performance of all IOTA strategies (Sayasneh *et al*, 2013a; Sayasneh *et al*, 2013b; Timmerman *et al*, 2010a; Van Holsbeke *et al*, 2012; Kaijser *et al*, 2013a; Timmerman *et al*, 2007; Timmerman *et al*, 2005; Timmerman *et al*, 2008; Ameye *et al*, 2012; Timmerman *et al*, 2010b; Van Holsbeke *et al*, 2009). On the other hand, validation studies of LR1, LR2 and Simple Rules performed outside IOTA studies reported similar results (Nunes *et al*, 2012; Alcazar *et al*, 2013; Hartman *et al*, 2012; Nunes *et al*, 2013), and there is now evidence that the IOTA strategies retain their performance in the hands of sonographers and

relatively inexperienced doctors (Sayasneh *et al*, 2013a; Sayasneh *et al*, 2013b; Nunes *et al*, 2012; Alcazar *et al*, 2013; Hartman *et al*, 2012; Nunes *et al*, 2013).

Our results showing that the IOTA methods and strategies have excellent ability to discriminate between benign and malignant adnexal masses and are superior to RMI in this regard are in line with other validation studies (Sayasneh *et al*, 2013a; Sayasneh *et al*, 2013b; Timmerman *et al*, 2010a; Van Holsbeke *et al*, 2012) and a recent systematic review (Kaijser *et al*, 2013b). The conclusion of the review was that an evidence-based approach to the preoperative characterisation of adnexal masses should incorporate the use of IOTA Simple Rules or LR2 instead of RMI, particularly in women of reproductive age (Kaijser *et al*, 2013b).

Our multicentre study demonstrated differences in test performance between centres for LR2, Simple Rules and RMI. However, in all centres, also in those with a low observed prevalence of malignancy, the sensitivity with regard to malignancy was much higher for the IOTA methods than for RMI. The overall discriminative capacity (AUC) for LR2 and RMI did not seem to be affected by cancer prevalence. However, our study highlighted important differences in calibration results for LR2, Simple Rules, and RMI. Type of centre appeared to contribute to these differences: oncology centres have a higher prevalence of malignant tumours, and suffered from underestimation of the predicted risk. This was also noticeable for Simple Rules and RMI even though these methods do not directly provide an estimated risk. For example, for RMI we can derive that the implicit average risk at an RMI value of 200 is 54%. This implies that at this (commonly used) cut-off, in some centres patients with a risk of malignancy of more than 50% may be classified as low risk.

We did not undertake a full meta-regression analysis (Van Houwelingen *et al*, 2002) to explain in detail the inter-centre differences in results for test performance and calibration as

this is beyond the scope of this paper. The most plausible explanations are differences in study populations (e.g. patients' age, Body Mass Index, or tumour mix), equipment, and examiners' use of the IOTA terms. However, the variation between centres of the observed (true risk) versus predicted risks for ovarian cancer revealed by this meta-analysis highlights that caution is needed when interpreting and using diagnostic test results (i.e. risks) for individual patient management within the context of a single centre. Future studies should explore the reasons for differences in diagnostic performance and calibration of diagnostic approaches between different centres, the final aim being to improve risk prediction for ovarian cancer.

Subjective assessment of grey scale and Doppler ultrasound findings by a very experienced ultrasound examiner has been suggested to be the preferred approach to characterise adnexal masses (Valentin *et al*, 2001). Unfortunately, most gynaecologists, radiologists, and sonographers have limited experience with the use of subjective assessment of ultrasound images to discriminate between benign and malignant adnexal masses. Because Simple Rules, Simple Descriptors and LR2 have been shown to perform very well in the hands of both sonographers and gynaecologists with limited ultrasound experience (Sayasneh *et al*, 2013a; Sayasneh *et al*, 2013b; Nunes *et al*, 2012; Alcazar *et al*, 2013; Hartman *et al*, 2012; Nunes *et al*, 2013), they could be used as first stage tests, and patients with inconclusive or equivocal results of the first stage test could be referred for subjective assessment by an experienced ultrasound examiner.

Each IOTA strategy has its own advantages and disadvantages. For example LR1 and LR2 give a continuous result (a risk estimate) for which the cut-off to diagnose malignancy can be varied depending on the context. The Simple Rules are easier to apply than LR1 and LR2, which require a computer or mobile application, but do not offer the flexibility of LR1 and

LR2. However, all of these approaches can be used to either classify all patients or classify a majority of patients while referring a subset of patients for further testing. Using Simple Descriptors as a first stage test offers no substantial advantages in test performance over the other IOTA strategies that we evaluated in this work. However, referring patients to expert examiners with masses in which the Simple Rules do not apply or with equivocal results of LR2 is advantageous as it leads to a reduction in the false positive rate whilst only minimally decreasing the sensitivity. In the current study as well as in IOTA phase 2 data, Simple Rules were inconclusive in 23% of patients whereas LR2 results were equivocal in 17-18% of the same patients (Timmerman *et al*, 2010a; Van Calster *et al*, 2012). In other validation studies of the Simple Rules fewer patients had inconclusive results, with reported percentages between 11% and 21% in different populations (Sayasneh *et al*, 2013b; Alcazar *et al*, 2013; Hartman *et al*, 2012; Fathallah *et al*, 2011).

The results of IOTA3 show that IOTA methods result in better discrimination of adnexal pathology prior to surgical treatment irrespective of the prevalence of malignant disease. Therefore, the application of IOTA risk models or rules provide a rational basis for referral of patients with a mass classified as malignant to specialist oncology services.

Acknowledgements

We are grateful to all study centres, investigators, and patients for their collaboration. BVC is a postdoctoral fellow of the Research Foundation – Flanders. LW is supported by a PhD grant of the Flanders' Agency for Innovation by Science and Technology (IWT Vlaanderen). This study is supported by FWO (project G049312N) and IWT (IWT-TBM 070706-IOTA3, the Swedish Medical Research Council (grants no. K2006-73X-11605-11-3); funds administered by Skåne University Hospital; Allmänna Sjukhusets i Malmö Stiftelse för bekämpande av cancer (the Malmö General Hospital Foundation for fighting against cancer); Landstingsfinansierad regional forskning and ALF-medel (i.e., two Swedish governmental grants from the region of Scania); Funds administered by Skåne University Hospital. TB is supported by Imperial Healthcare NHS Trust NIHR Biomedical Research Centre. The views expressed are those of the author(s) and not necessarily those of the NHS, the NIHR or the Department of Health.

Contributors

DT, ACT, BVC, IV, TB and LV conceived of the study and its design. ACT, DFi, CVH, DFr, LS, EE, AC, SG, RFr, FPGL, LV, JK and DT enrolled patients and acquired data. BVC, LW, CVH, JK, and DT were involved in data cleaning. BVC analysed the data, with support from LW. BVC, LW, LV, ACT, JK, TB, and DT were involved in data interpretation. ACT, JK, LW, BVC, LV, TB, and DT wrote the first draft of the manuscript which was then critically reviewed and revised by all other co-authors. All authors approved the final version of the manuscript for submission, and agree to be accountable for all aspects of the work relating to accuracy and integrity.

Conflict of interest

All authors state that they have no conflicts of interests to disclose.

References

Alcázar JL, Pascual MA, Olartecoechea B, Graupera B, Aubá M, Ajossa S, Hereter L, Julve R, Gastón B, Peddes C, Sedda F, Piras A, Saba L, Guerriero S (2013) IOTA simple rules for discriminating between benign and malignant adnexal masses: prospective external validation. *Ultrasound Obstet Gynecol* **42**(4):467-71

Ameye L, Timmerman D, Valentin L, Paladini D, Zhang J, Van Holsbeke C, Lissoni AA, Savelli L, Veldman J, Testa AC, Amant F, Van Huffel S, Bourne T (2012) Clinically oriented three-step strategy for assessment of adnexal pathology. *Ultrasound Obstet Gynecol* **40**(5):582-91

Bossuyt PM, Reitsma JB, Bruns DE, Gatsonis PP, Glasziou PP, Irwig LM, Lijmer JG, Moher D, Rennie D, De Vet HC (2003) Towards complete and accurate reporting of studies of diagnostic accuracy: The STARD Initiative. *Ann Intern Med* **138**:40-4

Bouwmeester W, Twisk JW, Kappen TH, van Klei WA, Moons KG, Vergouwe Y (2013) Prediction models for clustered data: comparison of a random intercept and standard regression model. *BMC Med Res Methodol* **13**:19

Cox DR (1958) Two further applications of a model for binary regression. *Biometrika* **45**:562-565

Deeks JJ (2001) Systematic reviews in health care: Systematic reviews of evaluations of diagnostic and screening tests. *BMJ* **323**(7305):157-62

Education, Practical Standards Committee, European Federation of Societies for Ultrasound in Medicine and Biology. Minimum training recommendations for the practice of medical ultrasound (2006) *Ultraschall in der Medizin* **27**(1):79-105

Fathallah K, Huchon C, Bats AS, Metzger U, Lefrere-Belda MA, Bensaid C, Lecuru F (2011) External validation of simple ultrasound rules of Timmerman on 122 ovarian tumors. *Gynecologie, Obstetrique & Fertilité* **39**(9):477-81

Hartman CA, Juliato CR, Sarian LO, Toledo MC, Jales RM, Morais SS, Pitta DD, Marussi EF, Derchain S (2012) Ultrasound criteria and CA 125 as predictive variables of ovarian cancer in women with adnexal tumors. *Ultrasound Obstet Gynecol* **40**(3):360-6

Heintz APM, Odicino F, Maisonneuve P, Beller U, Benedet JL, Creasman WT, Ngan HY, Pecorelli S (2003) Carcinoma of the Ovary: 25th Annual Report on the Results of Treatment in Gynecological Cancer. *Int J Gynaecol Obstet* **83**(suppl 1):S135-S137

Jacobs I, Oram D, Fairbanks J, Turner J, Frost C, Grudzinskas JG (1990) A risk of malignancy index incorporating CA 125, ultrasound and menopausal status for the accurate preoperative diagnosis of ovarian cancer. *B J Obstet Gynaecol* **97**(10):922-9

Kaijser J, Van Gorp T, Van Hoorde K, Van Holsbeke C, Sayasneh A, Vergote I, Bourne T, Timmerman D, Van Calster B (2013a) A comparison between an ultrasound based prediction model (LR2) and the Risk of Ovarian Malignancy Algorithm (ROMA) to assess the risk of malignancy in women with an adnexal mass. *Gynecol Oncol* **129**(2):377-83

Kaijser J, Sayasneh A, Van Hoorde K, Ghaem-Maghani S, Bourne T, Timmerman D, Van Calster B (2013b) Presurgical diagnosis of adnexal tumours using mathematical models and scoring systems: a systematic review and meta-analysis. *Human Reprod Update* doi:10.1093/humupd/dmt059 [Epub ahead of print]

Macaskill P, Gatsonis C, Deeks JJ, Harbord R, Takwoingi Y. Chapter 10: Analysing and Presenting Results. In: Deeks JJ, Bossuyt PM, Gatsonis C (editors), *Cochrane Handbook for Systematic Reviews of Diagnostic Test Accuracy* Version 1.0. The Cochrane Collaboration, 2010. Available from: <http://srdta.cochrane.org/>

Miller RW, Ueland FR (2012) Risk of malignancy in sonographically confirmed ovarian tumors. *Clin Obstet Gynecol* **55**(1):52-64

Moore RG, McMeekin DS, Brown AK, DiSilvestro P, Miller MC, Allard WJ, Gajewski W, Kurman R, Bast RC Jr, Skates SJ (2009) A novel multiple marker bioassay utilizing HE4 and CA125 for the prediction of ovarian cancer in patients with a pelvic mass. *Gynecol Oncol* **112**:40-6

Nunes N, Ambler G, Hoo WL, Naftalin J, Foo X, Widschwendter M, Jurkovic D (2013) A Prospective Validation of the IOTA Logistic Regression Models (LR1 and LR2) in Comparison to Subjective Pattern Recognition for the Diagnosis of Ovarian Cancer. *Int J Gynecol Cancer* **23**(9):1583-9

Nunes N, Yazbek J, Ambler G, Hoo W, Naftalin J, Jurkovic D (2012) A prospective evaluation of the IOTA Logistic Regression Model (LR2) for the diagnosis of ovarian cancer. *Ultrasound Obstet Gynecol* **40**(3):355-9

Riley RD, Dodd SR, Craig JV, Thompson JR, Williamson PR (2008) Meta-Analysis of Diagnostic Test Studies Using Individual Patient Data and Aggregate Data. *Stat Med* **27**(29):6111-36

Sayasneh A, Kaijser J, Preisler J, Johnson S, Stalder C, Husicka R, Guha S, Naji O, Abdallah Y, Raslan F, Drought A, Smith AA, Fotopoulou C, Ghaem-Maghami S, Van Calster B, Timmerman D, Bourne T (2013a) A multicenter prospective external validation

of the diagnostic performance of IOTA simple descriptors and rules to characterize ovarian masses. *Gynecol Oncol* **130**(1):140-6

Sayasneh A, Wynants L, Preisler J, Kaijser J, Johnson S, Stalder C, Husicka R, Abdallah Y, Raslan F, Drought A, Smith AA, Ghaem-Maghani S, Epstein E, Van Calster B, Timmerman D, Bourne T (2013b) Multicenter external validation of IOTA prediction models and RMI by operators with varied training. *Br J Cancer* **108**(12):2448-54

Sterne JA, White IR, Carlin JB, Spratt M, Royston P, Kenward MG, Wood AM, Carpenter JR (2009) Multiple imputation for missing data in epidemiological and clinical research: potential and pitfalls. *BMJ* **338**:b2393

Steyerberg EW (2009) *Clinical Prediction Models: A Practical Approach to Development, Validation, and Updating (Statistics for Biology and Health)*. Springer: New York

Timmerman D, Ameye L, Fischerova D, Epstein E, Melis GB, Guerriero S, Van Holsbeke C, Savelli L, Fruscio R, Lissoni AA, Testa AC, Veldman J, Vergote I, Van Huffel S, Bourne T, Valentin L (2010a) Simple ultrasound rules to distinguish between benign and malignant adnexal masses before surgery: prospective validation by IOTA group. *BMJ* **341**:c6839

Timmerman D, Testa AC, Bourne T, Ameye L, Jurkovic D, Van Holsbeke C, Paladini D, Van Calster B, Vergote I, Van Huffel S, Valentin L(2008) Simple ultrasound-based rules for the diagnosis of ovarian cancer. *Ultrasound Obstet Gynecol* **31**(6):681-90

Timmerman D, Testa AC, Bourne T, Ferrazzi E, Ameye L, Konstantinovic ML, Van Calster B, Collins WP, Vergote I, Van Huffel S, Valentin L; International Ovarian Tumor Analysis Group (2005) Logistic regression model to distinguish between the benign and

malignant adnexal mass before surgery: a multicenter study by the International Ovarian Tumor Analysis Group. *J Clin Oncol* **23**:8794–8801.

Timmerman D, Valentin L, Bourne TH, Collins WP, Verrelst H, Vergote I; International Ovarian Tumor Analysis (IOTA) Group.(2000) 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* **16**:500–505

Timmerman D, Van Calster B, Jurkovic D, Valentin L, Testa AC, Bernard JP, Van Holsbeke C, Van Huffel S, Vergote I, Bourne T (2007) Inclusion of CA-125 does not improve mathematical models developed to distinguish between benign and malignant adnexal tumors. *J Clin Oncol* **25**:4194–4200

Timmerman D, Van Calster B, Testa AC, Guerriero S, Fischerova D, Lissoni AA, Van Holsbeke C, Fruscio R, Czekierdowski A, Jurkovic D, Savelli L, Vergote I, Bourne T, Van Huffel S, Valentin L (2010b) Ovarian cancer prediction in adnexal masses using ultrasound-based logistic regression models: a temporal and external validation study by the IOTA group. *Ultrasound Obstet Gynecol* **36**(2):226-34

Valentin L, Hagen B, Tingulstad S, Eik-Nes S (2001) Comparison of ‘pattern recognition’ and logistic regression models for discrimination between benign and malignant pelvic masses: a prospective cross-validation. *Ultrasound Obstet Gynecol* **18**:357–65

Van Calster B, Timmerman D, Valentin L, McIndoe A, Ghaem-Maghami S, Testa AC, Vergote I, Bourne T (2012) Triaging women with ovarian masses for surgery: observational diagnostic study to compare RCOG guidelines with an International Ovarian Tumour Analysis (IOTA) group protocol. *BJOG* **119**(6):662-71

Van Calster B, Valentin L, Van Holsbeke C, Zhang J, Jurkovic D, Lissoni AA, Testa AC, Czekierdowski A, Fischerová D, Domali E, Van de Putte G, Vergote I, Van Huffel S, Bourne T, Timmerman D (2011) A novel approach to predict the likelihood of specific ovarian tumor pathology based on serum CA-125: a multicenter observational study. *Cancer Epidemiol Biomarkers Prev* **20**(11):2420-8

van Houwelingen HC, Arends LR, Stijnen T (2002) Advanced methods in meta-analysis: multivariate approach and meta-regression. *Stat Med* **21**(4):589-624.

Van Holsbeke C, Van Calster B, Bourne T, Ajossa S, Testa AC, Guerriero S, Fruscio R, Lissoni AA, Czekierdowski A, Savelli L, Van Huffel S, Valentin L, Timmerman D (2012) External validation of diagnostic models to estimate the risk of malignancy in adnexal masses. *Clin Cancer Res* **18**(3):815-25

Van Holsbeke C, Van Calster B, Testa AC, Domali E, Lu C, Van Huffel S, Valentin L, Timmerman D (2009) Prospective internal validation of mathematical models to predict malignancy in adnexal masses: results from the international ovarian tumor analysis study. *Clin Cancer Res* **15**:684–691

Van Klaveren D, Steyerberg EW, Perel P, Vergouwe Y (2014) Assessing discriminative ability of risk models in clustered data. *BMC Med Res Methodol* **14**(1):5. doi: 10.1186/1471-2288-14-5

Verleye L, Vergote I, van der Zee AG (2010) Patterns of care in surgery for ovarian cancer in Europe. *Eur J Surg Oncol* **36** Suppl 1:S108-14

Woo YL, Kyrgiou M, Bryant A, Everett T, Dickinson HO (2012) Centralisation of services for gynaecological cancers - a Cochrane systematic review. *Gynecol Oncol* **126**(2):286-90

Table 1. Description of the International Ovarian Tumour Analysis (IOTA) methods evaluated in the IOTA study phase 3.

IOTA method	Variables or features
Logistic regression model-1 (LR1) ¹⁰ (risks \geq 10% indicate malignancy)	(1) personal history of ovarian cancer (yes, 1; no, 0), (2) current use of hormonal therapy (yes, 1; no, 0), (3) age of the patient (in years), (4) maximum diameter of lesion (in mm), (5) tender mass at examination (yes, 1; no, 0), (6) ascites (yes, 1; no, 0), (7) blood flow in papillary projection (yes, 1; no, 0), (8) purely solid tumour, (9) maximum diameter of the largest solid component (in mm, but with no increase >50 mm), (10) irregular internal cyst walls (yes, 1; no, 0), (11) acoustic shadows (yes, 1; no, 0), and (12) colour flow score (1–4, where 1 is no flow and 4 is maximum flow) The mathematical formula is presented in Supplementary Appendix
Logistic regression model-2 (LR2) ¹⁰ , (risks \geq 10% indicate malignancy)	(1) ascites (yes, 1; no, 0), (2) blood flow in papillary projection (yes, 1; no, 0), (3) maximum diameter of the largest solid component (in mm, but with no increase >50 mm), (4) irregular internal cyst walls (yes, 1; no, 0), (5) acoustic shadows (yes, 1; no, 0), (6) age of the patient (in years) The mathematical formula is presented in Supplementary Appendix
IOTA Simple Rules (SR) ^{11, a}	Benign features: Unilocular tumour (B1), Largest diameter of largest solid component <7mm (B2), Acoustic shadows (B3), Smooth multilocular tumour with largest diameter <100 mm (B4), No intratumoral blood flow at colour or power Doppler (B5). Malignant features: Irregular solid tumour (M1), Ascites (M2), At least 4 papillary projections (M3), Irregular multilocular solid tumour with largest diameter \geq 100 mm (M4), Very strong intratumoral blood flow at colour or power Doppler (M5).
IOTA Simple Descriptors (SD) ^{12, b}	Benign Descriptors: Unilocular tumour with ground glass echogenicity in a premenopausal woman; Unilocular tumour with mixed echogenicity and acoustic shadows in a premenopausal woman; Unilocular anechoic tumour with regular walls and maximum diameter of lesion < 10 cm; Remaining unilocular tumours with regular walls Malignant Descriptor: Tumour with ascites and at least moderate colour Doppler blood flow in a postmenopausal woman; Age > 50 years and CA 125 > 100 U/mL.

^a A mass is classified as malignant if at least one M-feature and none of the B-features are present and vice versa. If no B or M features are present, or if both B and M features are

present, then the rules are considered inconclusive (unclassifiable mass), and a second stage test should be used in the unclassifiable tumours.

^b A mass classified as malignant if at least one malignant descriptor and none of the benign descriptors are present and vice versa. If no benign or malignant descriptors are present, or if both benign and malignant descriptors are present, then the descriptors are inconclusive (unclassifiable mass), and a second stage test should be used in the unclassifiable tumours.

Table 2. Overview of tumour types.

Pathology	Frequency (n)	Percent (%)
Benign masses		
Endometrioma	344	14.3
Serous cystadenoma	259	10.8
Teratoma	231	9.6
Mucinous cystadenoma	183	7.6
Fibroma	130	5.4
Simple cyst or parasalpingeal cyst	106	4.4
Rare benign ^a	48	2.0
Hydrosalpinx or salpingitis	47	2.0
Functional cyst	40	1.7
Peritoneal pseudocyst	18	0.8
Abscess	17	0.7
Malignant masses		
Primary invasive stage I	128	5.3
Primary invasive stage II	47	2.0
Primary invasive stage III	397	16.5
Primary invasive stage IV	61	2.5
Borderline stage I	135	5.6
Borderline stage II	6	0.3
Borderline stage III	12	0.5
Rare primary invasive ^b	68	2.8
Metastatic	126	5.2
Total	2403	100

^a For example Brenner tumour or struma ovarii

^b For example dysgerminoma, granulosa cell tumour, yolk sac tumour, or malignant teratoma

Table 3. Results with regard to the variables included in LR1, LR2 and the Simple Rules in benign and malignant adnexal masses

Variable	Statistics	Benign (1423, 59%)	Malignant (980, 41%)
Variables in LR1 and LR2			
Age (years)	Median (IQR)	44 (33-56)	57 (46-66)
Largest diameter of lesion (mm)	Median (IQR)	64 (47-90)	86 (56-126)
Solid components	N (%)	472 (33%)	915 (93%)
Largest diameter of solid component if present (mm)	Median (IQR)	28 (13-54)	59 (37-87)
Colour score (1-4)			
Colour score 1	N (%)	574 (40)	32 (3)
Colour score 2	N (%)	563 (40)	199 (20)
Colour score 3	N (%)	239 (17)	442 (45)
Colour score 4	N (%)	47 (3)	307 (31)
Ascites	N (%)	18 (1%)	322 (33%)
Papillations with detectable blood flow	N (%)	55 (4%)	160 (16%)
Irregular cyst walls	N (%)	385 (27%)	572 (58%)
Acoustic shadows	N (%)	265 (19%)	34 (3%)
Tender mass at ultrasound examination	N (%)	233 (16%)	111 (11%)
Current use of hormonal therapy	N (%)	153 (11%)	54 (6%)
Personal history of ovarian cancer	N (%)	14 (1%)	30 (3%)
Solid tumour	N (%)	154 (11%)	473 (48%)
Variables in the simple rules			
Benign ultrasound features in the simple rules			
Unilocular tumour (B1)	N (%)	595 (42%)	5 (0.5%)
Largest diameter of largest solid	N (%)	40 (3%)	2 (0.2%)

component <7mm (B2)

Variable	Statistics	Benign (1423, 59%)	Malignant (980, 41%)
Acoustic shadows (B3)	N (%)	265 (19%)	34 (3%)
Smooth multilocular tumour with largest diameter <100 mm (B4)	N (%)	224 (16%)	13 (1%)
No intratumoral blood flow at colour or power Doppler (B5)	N (%)	574 (40%)	32 (3%)
Malignant ultrasound features in the simple rules			
Irregular solid tumour (M1)	N (%)	16 (1%)	189 (19%)
Ascites (M2)	N (%)	18 (1%)	322 (33%)
At least 4 papillary projections (M3)	N (%)	27 (2%)	91 (9%)
Irregular multilocular solid tumour with largest diameter \geq 100 mm (M4)	N (%)	40 (3%)	153 (16%)
Very strong intratumoral blood flow at colour or power Doppler (M5)	N (%)	47 (3%)	307 (31%)

IQR: interquartile range; SD: standard deviation

Table 4. Test performance of the International Ovarian Tumour Analysis (IOTA) diagnostic strategies, subjective assessment and Risk of Malignancy Index (RMI) when using a meta-analysis approach on centre-specific data.

Approach	AUC (95%CI)	Sens, % (95% CI)	Spec, % (95% CI)	LR+	LR-	DOR (95% CI)
One-step strategies						
LR1	0.930 (0.917-0.942)	93.7 (91.4-95.4)	77.6 (70.9-83.0)	4.17	0.08	40.8 (30.0-55.4)
LR2	0.918 (0.905-0.930)	90.2 (86.9-92.8)	78.9 (73.2-83.7)	4.28	0.12	31.2 (23.1-42.2)
SA	0.914 (0.886-0.936)	92.5 (89.4-94.8)	87.7 (83.2-91.2)	7.53	0.09	72.9 (49.8-107)
RMI	0.875 (0.853, 0.894)	67.1 (61.4-72.4)	90.6 (87.3-93.1)	7.15	0.36	17.5 (13.1-23.4)
SRMal	N/A	95.3 (93.1-96.9)	74.1 (67.7-79.7)	3.68	0.06	49.1 (34.9-69.0)
Two step strategies						
SR+SA	N/A	91.8 (89.1-93.9)	89.0 (85.2-92.0)	8.38	0.09	75.7 (55.6-103)
LR2+SA	N/A	92.3 (89.5-94.5)	84.8 (80.4-88.3)	6.06	0.09	58.7 (43.4-79.4)
SD+SRMal	N/A	95.7 (93.5-97.1)	73.6 (66.7-79.5)	3.62	0.06	50.5 (35.7-71.6)
SD+LR2	N/A	91.1 (88.1-93.5)	78.1 (72.4-82.9)	4.17	0.11	32.8 (24.6-43.7)
SD+SA	N/A	93.0 (90.0-95.1)	86.5 (81.8-90.1)	6.88	0.08	68.5 (47.7-98.3)
						cont

Table 4 Cont.						
Approach	AUC (95%CI)	Sens, % (95% CI)	Spec, % (95% CI)	LR+	LR-	DOR (95% CI)
Three-step strategies						
SD+SR+SA	N/A	92.5 (89.6-94.6)	87.6 (83.5-90.7)	7.44	0.09	70.7 (51.7-96.5)
SD+LR2+ SA	N/A	93.1 (90.5-95.0)	83.7 (79.2-87.4)	5.71	0.08	57.6 (42.3-78.6)

AUC, area under the receiver-operating characteristics curve; Sens, sensitivity; Spec, specificity; LR+, positive likelihood ratio; LR-, negative likelihood ratio; DOR: diagnostic odds ratio; CI, confidence interval; N/A, not applicable. LR1, logistic regression model 1; LR2, logistic regression model 2; SA, subjective assessment; RMI-1, risk of malignancy index-1; SR, simple rules; SD, simple descriptors; SRMal, Simple Rules as a first stage test with all tumours in which Simple Rules are inconclusive being classified as malignant; SR+SA, Simple Rules as a first stage test and subjective assessment for tumours in which the Simple Rules are inconclusive; LR2+SA, LR2 as a first stage test and subjective assessment for tumours in which LR2 yields a predicted risk of malignancy of $\geq 5\%$ but $< 25\%$; SD+SRMal, Simple Descriptors as a first stage test, Simple Rules for tumours unclassifiable by the descriptors with all tumours in which Simple Rules are inconclusive being classified as malignant; SD+LR2, Simple Descriptors as a first stage test and LR2 for those tumours in which the descriptors are not applicable; SD+SA, Simple Descriptors as a first stage test and subjective assessment for those tumours in which the Simple Descriptors are not applicable; SD+SR+SA: Simple Descriptors as a first stage test, Simple Rules for tumours in which the descriptors are not applicable, and subjective assessment for masses in which the Simple

Rules are inconclusive; SD+LR2+SA, Simple Descriptors as a first stage test, LR2 for tumours in which the descriptors are not applicable, and subjective assessment for masses in which LR2 yields a predicted risk of $\geq 5\%$ but $< 25\%$.

Table 5. Test performance of LR2, Simple Rules and Risk of Malignancy Index (RMI) in pre- and postmenopausal patients using a meta-analysis approach on centre-specific data.

Diagnostic method	AUC (95%CI)	Sens, % (95%CI)	Spec, % (95%CI)
<i>Premenopausal patients</i>			
LR2	0.908 (0.886-0.926)	85 (78-90)	82 (77-87)
SRMal		95 (91-97)	77 (70-83)
SR+SA		92 (86-95)	91 (87-94)
RMI	0.867 (0.837-0.892)	53 (45-61)	94 (92-96)
<i>Postmenopausal patients</i>			
LR2	0.897 (0.872-0.917)	94 (92-96)	65 (58-71)
SRMal		96 (93-97)	66 (59-73)
SR+SA		93 (90-95)	83 (78-87)
RMI	0.850 (0.805-0.887)	78 (72-83)	81 (76-85)

AUC: area under the receiver-operating characteristics curve; Sens: sensitivity; Spec: specificity; CI: confidence interval. LR2: logistic regression model 2; SRMal: a one-step strategy using the IOTA Simple Rules as a first stage test and classifying tumours where the simple rules yield an inconclusive result as malignant; SR+SA: a two-stage strategy using the IOTA Simple Rules as a first stage test and using subjective assessment for those tumours where the simple rules yield an inconclusive result; RMI-1: risk of malignancy index-1

Titles and legends to figures

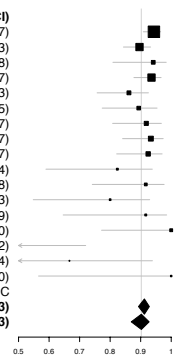
Figure 1. The sensitivity (Sens) and specificity (Spec) for LR2 (A), Risk of Malignancy Index (B) and a two-stage strategy using Simple Rules as a first stage test and using subjective assessment for tumours in which the Simple Rules are inconclusive (C) per contributing centre and for all centres combined using a meta-analysis approach and pooled data.

NC, not computable; CI, confidence interval. Numbers in brackets denote the prevalence (%) of malignant masses in each centre. **Oncology centres** were: University Hospitals Leuven, Belgium (LBE); Università Cattolica del Sacro Cuore, Rome, Italy (RIT); Ospedale San Gerardo, Monza, Italy (OIT); General Faculty Hospital, Prague, Czech Republic (PCR); Istituto Europeo di Oncologia, Milan, Italy (CIT); Medical University Lublin, Poland (LPO); Karolinska University Hospital, Stockholm, Sweden (SSW); Skåne University Hospital Lund, Sweden (LSW); Università degli Studi di Udine, Italy (UDI); Istituto Nazionale dei Tumori, Naples, Italy (GIT); University of Bologna, Italy (BIT). **Non-oncology centres** were: Skåne University Hospital Malmö, Sweden (MSW), Ziekenhuis Oost-Limburg, Genk, Belgium (GBE); Ospedale San Giovanni di Dio, Cagliari, Italy (SIT); DCS Sacco University of Milan, Italy (MIT); Università degli Studi di Napoli, Naples, Italy (NIT); Institut Universitari Dexeus, Barcelona, Spain (BSP); Ospedale dei Bambini Vittore Buzzi, Milan, Italy (FIT).

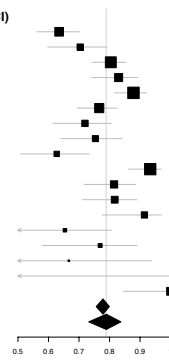
Figure 2. Centre-specific calibration curves (Cox logistic recalibration) for Risk of Malignancy Index (RMI), LR2, and proportions for Simple Rules. Vertical lines crossing the x-axis for RMI (200) and LR2 (0.1, i.e. 10% risk) represent the original cut-off to define malignant disease. Scatterplots above and below the calibration curves for RMI and LR2 represent the distribution of predicted risks for LR2 and values for RMI for benign and malignant tumours, respectively. The x-axis for RMI is limited to 1000. **Oncology centres**

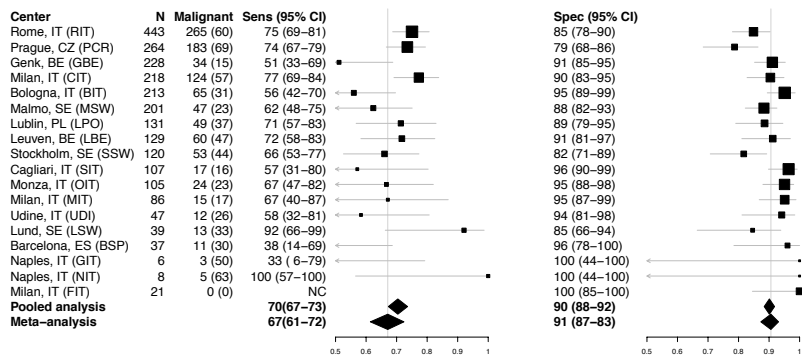
were: Rome, Prague, Milan, Leuven, Stockholm, Bologna and Lublin. **Non-oncology centres** were: Malmö and Genk.

Center	N	Malignant	Sens (95% CI)
Rome, IT (RIT)	443	265 (60)	94 (91-97)
Prague, CZ (PCR)	264	183 (69)	90 (84-93)
Genk, BE (GBE)	228	34 (15)	94 (81-98)
Milan, IT (CIT)	218	124 (57)	94 (88-97)
Bologna, IT (BIT)	213	65 (31)	86 (76-93)
Malmö, SE (MSW)	201	47 (23)	89 (77-95)
Lublin, PL (LPO)	131	49 (37)	92 (81-97)
Leuven, BE (LBE)	129	60 (47)	93 (84-97)
Stockholm, SE (SSW)	120	53 (44)	92 (82-97)
Cagliari, IT (SIT)	107	17 (16)	82 (59-94)
Monza, IT (OIT)	105	24 (23)	92 (74-98)
Milan, IT (MIT)	86	15 (17)	80 (55-93)
Udine, IT (UDI)	47	12 (26)	92 (65-99)
Lund, SE (LSW)	39	13 (33)	100 (77-100)
Barcelona, ES (BSP)	37	11 (30)	45 (21-72)
Naples, IT (GIT)	6	3 (50)	67 (21-94)
Naples, IT (NIT)	8	5 (63)	100 (57-100)
Milan, IT (FIT)	21	0 (0)	NC
Pooled analysis			91 (89-93)
Meta-analysis			90 (87-93)

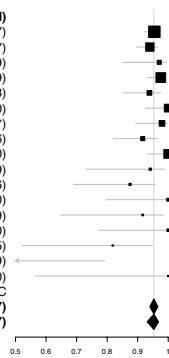


Spec (95% CI)
63 (56-70)
70 (60-79)
80 (74-85)
83 (74-89)
88 (82-92)
77 (69-83)
72 (61-81)
75 (64-84)
63 (51-73)
93 (86-97)
81 (72-88)
82 (71-89)
91 (78-97)
65 (46-81)
77 (58-89)
67 (21-94)
100 (44-100)
100 (85-100)
78 (76-80)
79 (73-84)

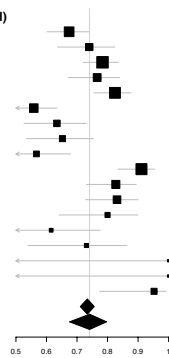


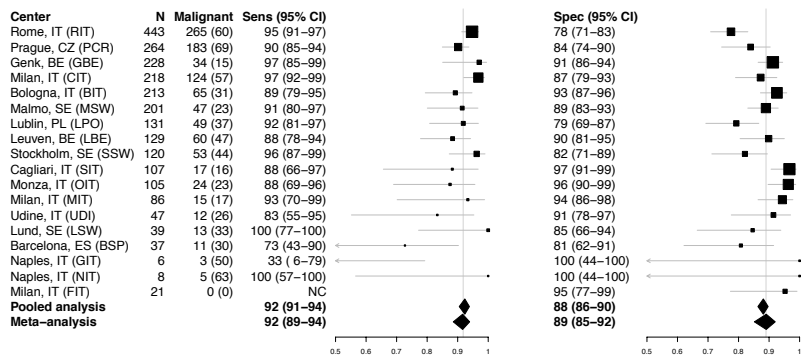


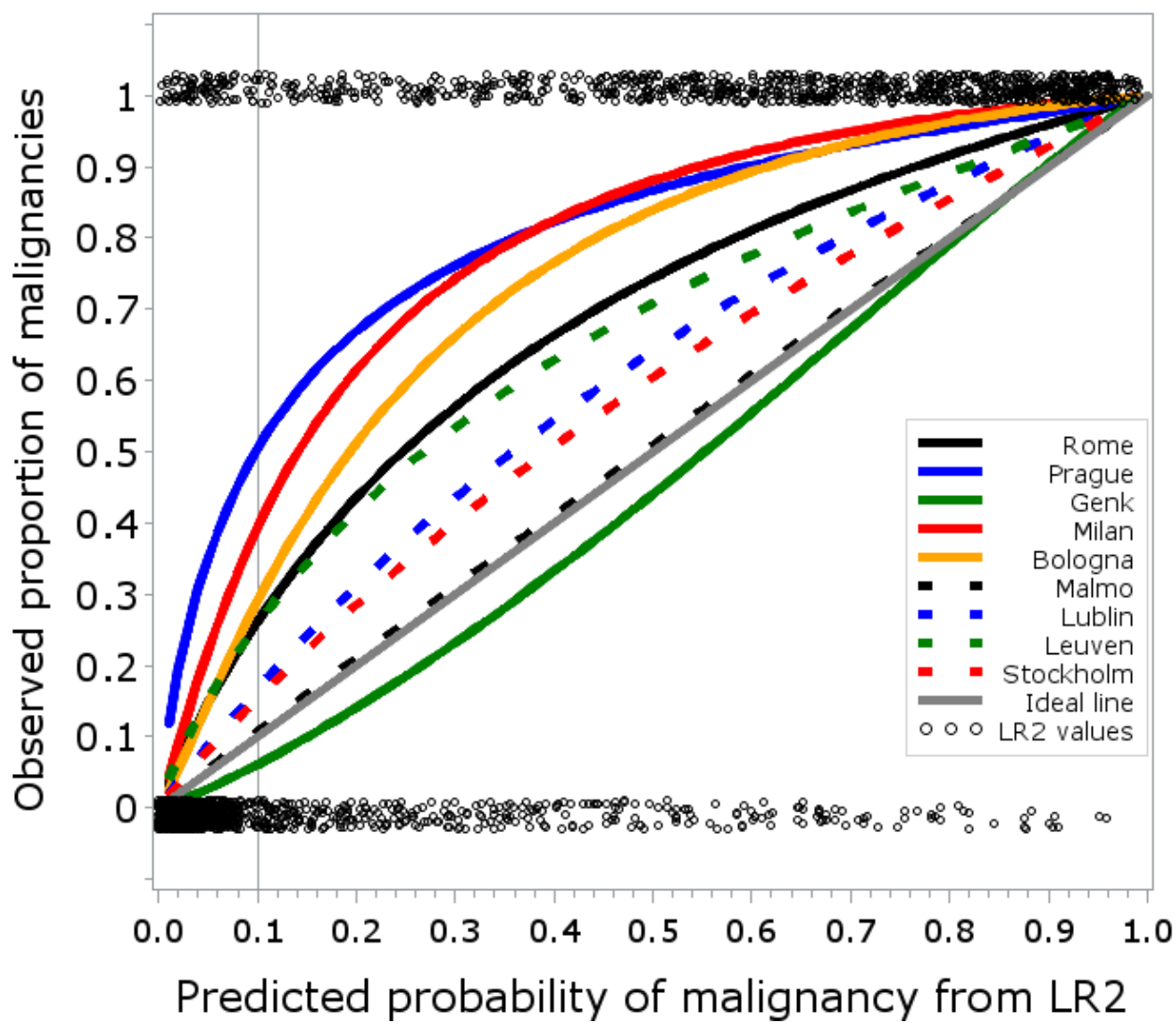
Center	N	Malignant	Sens (95% CI)
Rome, IT (RIT)	443	265 (60)	95 (92-97)
Prague, CZ (PCR)	264	183 (69)	94 (90-97)
Genk, BE (GBE)	228	34 (15)	97 (85-99)
Milan, IT (CIT)	218	124 (57)	98 (93-99)
Bologna, IT (BIT)	213	65 (31)	94 (85-98)
Malmö, SE (MSW)	201	47 (23)	100 (92-100)
Lublin, PL (LPO)	131	49 (37)	98 (89-100)
Leuven, BE (LBE)	129	60 (47)	92 (82-96)
Stockholm, SE (SSW)	120	53 (44)	100 (93-100)
Cagliari, IT (SIT)	107	17 (16)	94 (73-99)
Monza, IT (OIT)	105	24 (23)	88 (69-96)
Milan, IT (MIT)	86	15 (17)	100 (80-100)
Udine, IT (UDI)	47	12 (26)	92 (65-99)
Lund, SE (LSW)	39	13 (33)	100 (77-100)
Barcelona, ES (BSP)	37	11 (30)	82 (52-95)
Naples, IT (GIT)	6	3 (50)	33 (6-79)
Naples, IT (NIT)	8	5 (63)	100 (57-100)
Milan, IT (FIT)	21	0 (0)	NC
Pooled data			95 (94-97)
Meta-analysis			95 (93-97)

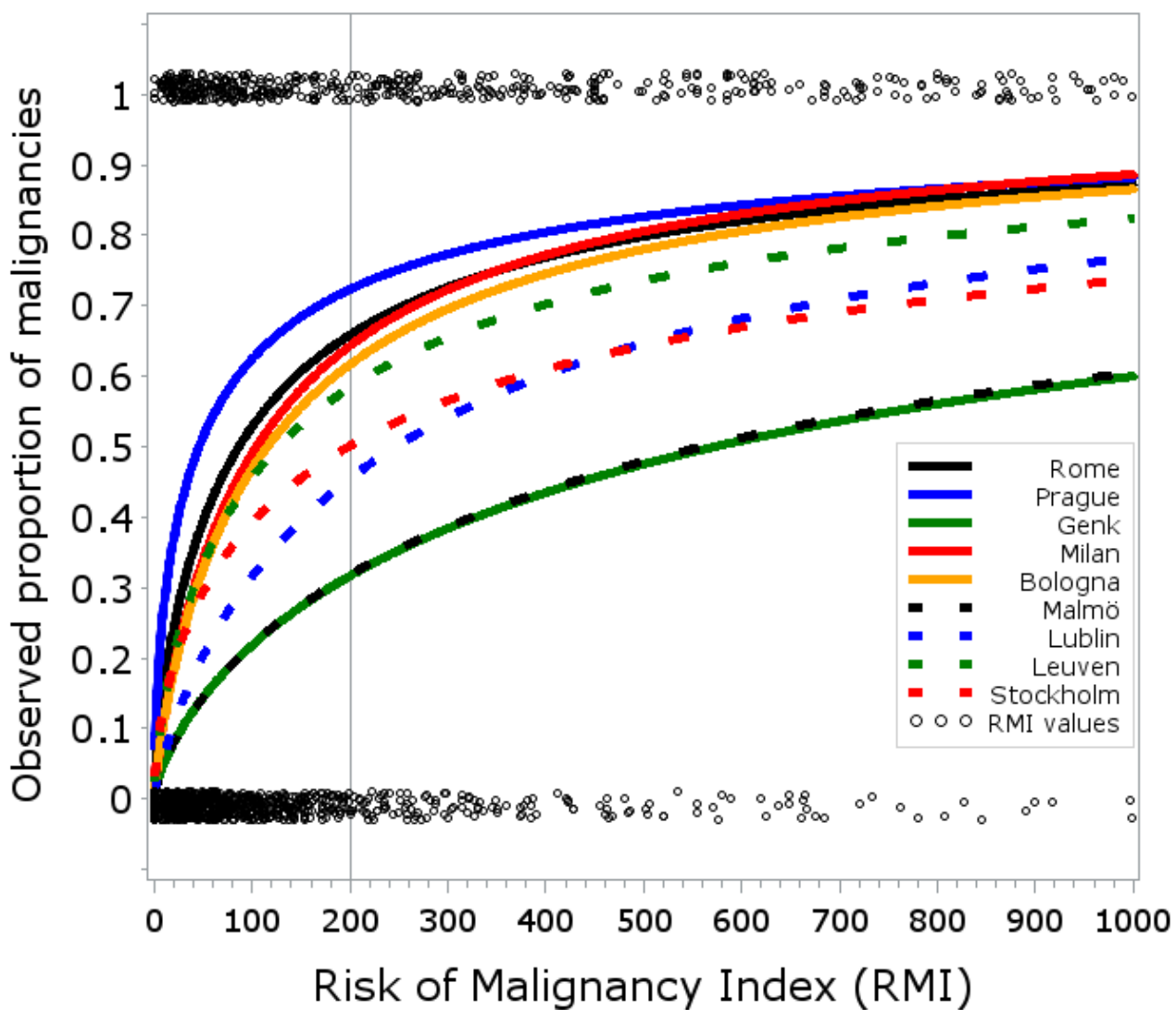


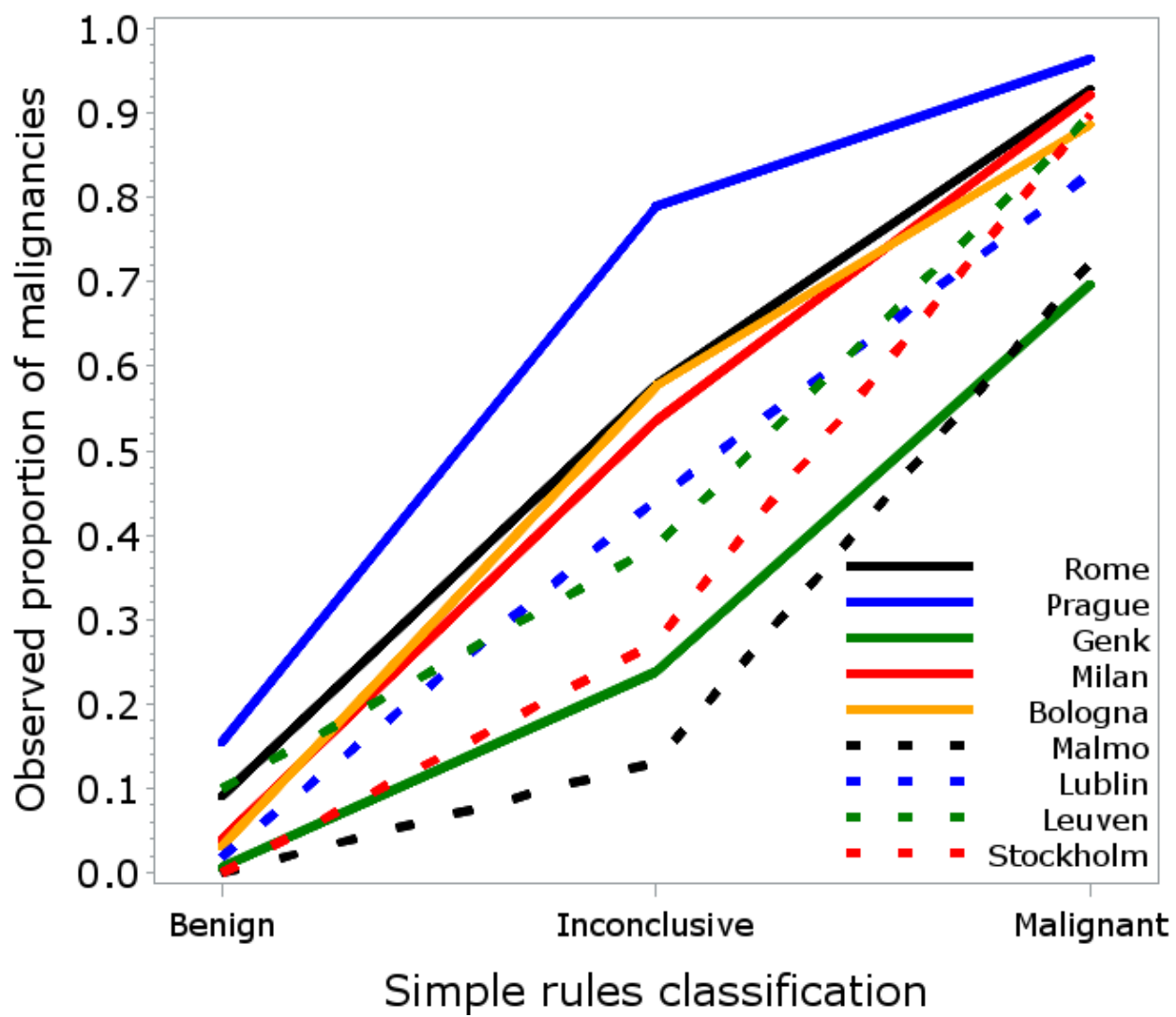
Spec (95% CI)
67 (60-74)
74 (64-82)
78 (72-84)
77 (67-84)
82 (76-88)
56 (48-63)
63 (53-73)
65 (53-75)
57 (45-68)
91 (83-95)
83 (73-89)
83 (73-90)
80 (64-90)
62 (43-78)
73 (54-86)
100 (44-100)
100 (44-100)
95 (77-99)
73 (71-76)
74 (68-80)











Appendix

Recruitment Centres in IOTA phase 3:

Participating in IOTA phase 1:

University Hospitals Leuven, Belgium (LBE); Università degli Studi di Napoli, Naples, Italy (NIT); Ospedale San Gerardo, Monza, Italy (OIT); DSC L. Sacco, University of Milan, Italy (MIT); Università Cattolica del Sacro Cuore, Rome, Italy (RIT); Skåne University Hospital Malmö, Sweden (MSW);

Participating in IOTA phase 1b:

University Hospitals Leuven, Belgium (LBE); Università Cattolica del Sacro Cuore, Rome, Italy (RIT); Skåne University Hospital Malmö, Sweden (MSW).

Participating in IOTA phase 2:

University Hospitals Leuven, Belgium (LBE); Ospedale San Gerardo, Monza, Italy (OIT); Ziekenhuis Oost-Limburg, Genk, Belgium (GBE); Medical University Lublin, Poland (LPO); Ospedale San Giovanni di Dio, Cagliari, Italy (SIT); Skåne University Hospital Malmö, Sweden (MSW); University of Bologna, Italy (BIT); Università Cattolica del Sacro Cuore, Rome, Italy (RIT); DSC L. Sacco, University of Milan, Italy (MIT); General Faculty Hospital, Prague, Czech Republic (PCR); Università degli Studi di Napoli, Naples, Italy (NIT); Istituto Europeo di Oncologia, Milan, Italy (CIT); Skåne University Hospital Lund, Sweden (LSW); Università degli Studi di Udine, Italy (UDI); Istituto Nazionale dei Tumori, Naples, Italy (GIT).

Participating only in IOTA phase 3:

Karolinska University Hospital, Stockholm, Sweden (SSW); Ospedale dei Bambini Vittore Buzzi, Milan, Italy (FIT); Institut Universitari Dexeus, Barcelona, Spain (BSP).

Oncology Centres in IOTA 3 (n=11)

University Hospitals Leuven, Belgium (LBE); Università Cattolica del Sacro Cuore, Rome, Italy (RIT); Ospedale San Gerardo, Monza, Italy (OIT); General Faculty Hospital, Prague, Czech Republic (PCR); Istituto Europeo di Oncologia, Milan, Italy (CIT); Medical University Lublin, Poland (LPO); Karolinska University Hospital, Stockholm, Sweden (SSW); Skåne University Hospital Lund, Sweden (LSW); Università degli Studi di Udine, Italy (UDI); Istituto Nazionale dei Tumori, Naples, Italy (GIT) University of Bologna, Italy (BIT)

Non-Oncology Centres in IOTA3 (n=7)

Skåne University Hospital Malmö, Sweden (MSW); Ziekenhuis Oost-Limburg, Genk, Belgium (GBE); University of Bologna, Italy (BIT); Ospedale San Giovanni di Dio, Cagliari, Italy (SIT); DSC L. Sacco, University of Milan, Italy (MIT); Università degli Studi di Napoli, Naples, Italy (NIT); Institut Universitari Dexeus, Barcelona, Spain (BSP); Ospedale dei Bambini Vittore Buzzi, Milan, Italy (FIT)

Mathematical formulas for IOTA LR1 and LR2:

LR 2: For the logistic regression model LR2 the estimated probability of malignancy for a patient with an adnexal tumour was equal to $y = 1/(1 + e^{-z})$, where $z = -5.3718 + 0.0354 (1) + 1.6159 (2) + 1.1768 (3) + 0.0697 (4) + 0.9586 (5) - 2.9486 (6)$, and e is the mathematical constant and base value of natural logarithms.

1) age of the patient (in years), (2) the presence of ascites (yes = 1, no = 0), (3) the presence of blood flow within a solid papillary projection (yes = 1, no = 0), (4) maximal diameter of the solid component (expressed in millimeters, but with no increase > 50 mm), (5) irregular internal cyst walls (yes = 1, no = 0), and (6) the presence of acoustic shadows (yes = 1, no = 0).

LR1: For the logistic regression model LR1 the estimated probability of malignancy for a patient with an adnexal tumor was equal to $y = 1/(1 + e^{-z})$, where $z = -6.7468 + 1.5985 (1) - 0.9983 (2) + 0.0326 (3) + 0.00841 (4) - 0.8577 (5) + 1.5513 (6) + 1.1737 (7) + 0.9281 (8) + 0.0496 (9) + 1.1421 (10) - 2.3550 (11) + 0.4916 (12)$, and e is the mathematical constant and base value of natural logarithms.

1) personal history of ovarian cancer (yes = 1, no = 0), (2) current hormonal therapy (yes = 1, no = 0), (3) age of the patient (in years), (4), maximum diameter of the lesion (in millimeters), (5) the presence of pain during the examination (yes = 1, no = 0), (6) the presence of ascites (yes = 1, no = 0), (7) the presence of blood flow within a solid papillary projection (yes = 1, no = 0), (8) the presence of a purely solid tumor (yes = 1, no = 0), (9) maximal diameter of the solid component (expressed in millimeters, but with no increase > 50 mm), (10) irregular internal cyst walls (yes = 1, no = 0), (11) the presence of acoustic shadows (yes = 1, no = 0), and (12) the color score (1, 2, 3, or 4).

Multiple imputation of missing values for serum CA-125

We used the IOTA data from phases 1, 1b, 2, and 3 to impute missing values. In these data, the serum CA-125 level was missing in 31% of the women. The value was more often missing in women with a benign tumour than in women with a malignant tumour. It is highly likely that missing values have occurred for two main reasons. Firstly, some centres were less committed than others to measure CA-125 due to differing management practices. Secondly, investigators sometimes decided not to measure CA-125 based on the general clinical picture of the patient and on the appearance of the tumour on ultrasound. We used the approach of ‘multiple imputation’ to deal with missing values in the analysis: the missing CA-125 values were estimated (i.e. imputed) more than once to acknowledge that we do not know what exact value we would have observed if the CA-125 level were available.¹ To estimate the missing values, we used predictive mean matching regression² using tumour histology and other variables that were related to either the level of CA-125 itself or to the unavailability of CA-125 (i.e. a binary indicator indicating for each woman whether CA-125 was missing or not). This was repeated 100 times to generate 100 completed data sets. The RMI was evaluated on all 100 completed datasets and the results were combined using the standard Rubin’s rules.

References

1. Sterne JAC, White IR, Carlin JB, Spratt M, Royston P, Kenward MG, et al. Multiple imputation for missing data in epidemiological and clinical research: potential and pitfalls. *BMJ* 2009;338:b2393.

2. Schenker N, Taylor JMG. Partially parametric techniques for multiple imputation. *Comput Stat Data Anal* 1996;22:425-46.

Other references regarding imputation: Ali *et al* *BJC* 2011; Clark *et al* *BJC* 2001; Donders *et al* *JCE* 2006; Little *et al* *NEJM* 2012

Subgroup analysis: discrimination of benign versus subtypes of malignant disease

(AUC, area under the receiver operating characteristic curve; sens, sensitivity; spec, specificity; CI, confidence interval; LR2, logistic regression model2; RMI, risk of malignancy index; SR, simple rules; SR+mal, simple rules as first stage test with all cases in which the simple rules cannot be applied classified as malignant; SR+SA, simple rules as first stage test and all cases in which the simple rules cannot be applied classified by subjective assessment)

Table S1. Benign (n=1423) versus borderline (n=153)

	AUC (95% CI)	Sens (95% CI)	Spec (95% CI)	N false negatives
LR2	0.819 (0.775-0.856)	73.1 (64.0-80.6)	78.5 (72.0-83.8)	45
RMI	0.694 (0.624-0.756)	29.6 (21.2-39.7)	90.6 (87.1-93.2)	110
SR+Mal	/	87.5 (79.3-92.8)	74.2 (66.5-80.7)	20
SR+SA	/	79.5 (70.8-86.1)	89.3 (84.7-92.7)	32

Table S2. Benign (n=1423) versus stage I invasive (n=128)

	AUC (95% CI)	Sens (95% CI)	Spec (95% CI)	N false negatives
LR2	0.908 (0.882-0.929)	88.7 (88.0-93.9)	78.8 (72.5-83.9)	13
RMI	0.810 (0.744-0.862)	58.2 (47.1-68.6)	90.5 (87.1-93.1)	53
SR+Mal	/	92.4 (84.6-96.4)	74.1 (67.1-80.0)	10
SR+SA	/	89.5 (81.3-94.4)	89.3 (85.2-92.3)	13

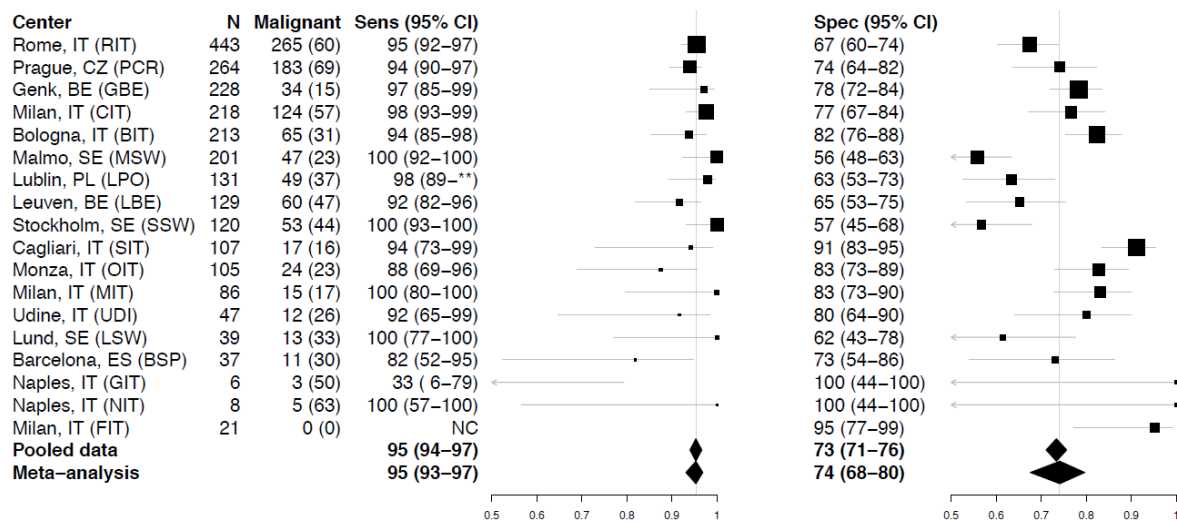
Table S3. Benign (n=1423) versus stage II-stage IV invasive (n=505)

	AUC (95% CI)	Sens (95% CI)	Spec (95% CI)	N false negatives
LR2	0.963 (0.948-0.973)	97.2 (94.1-98.7)	78.5 (72.1-83.8)	12
RMI	0.955 (0.938-0.968)	91.5 (87.4-94.4)	90.5 (87.1-93.1)	38
SR+Mal	/	99.2 (97.6-99.8)	74.3 (66.6-80.8)	4
SR+SA	/	98.6 (94.4-99.7)	89.4 (84.4-92.7)	7

Table S4. Benign (n=1423) versus metastatic (n=126)

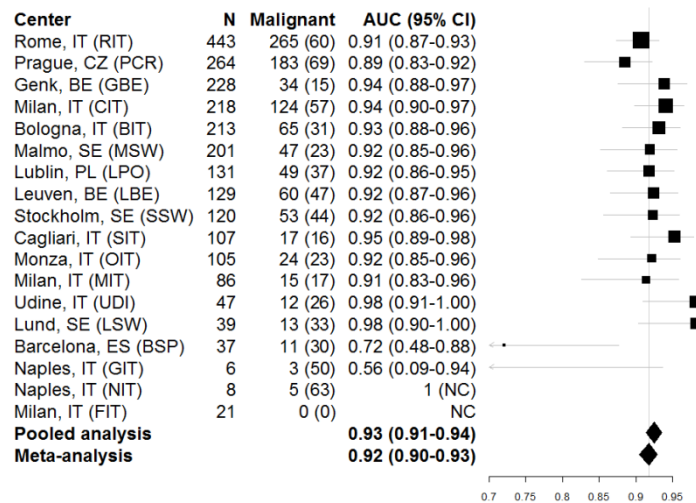
	AUC (95% CI)	Sens (95% CI)	Spec (95% CI)	N false negatives
LR2	0.925 (0.898-0.946)	93.3 (86.1-96.9)	78.5 (72.0-83.8)	9
RMI	0.842 (0.780-0.890)	65.9 (52.9-76.9)	90.7 (87.3-93.3)	47
SR+Mal	/	93.5 (81.7-97.9)	74.3 (66.6-80.8)	9
SR+SA	/	93.9 (75.4-98.7)	89.3 (84.7-92.7)	12

Figure S1. The sensitivity and specificity with 95% CI of using Simple Rules (classifying all tumours where the Simple Rules yield an inconclusive result as malignant) per contributing centre and for all centres combined using a meta-analysis approach and pooled data.



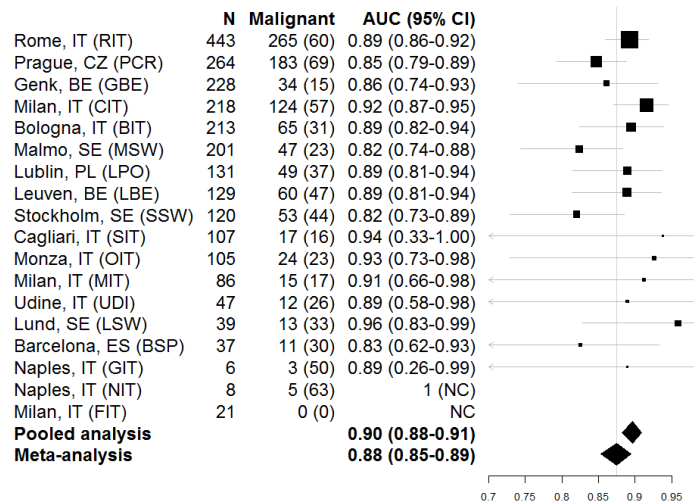
NC, not computable; CI, confidence interval. Numbers in brackets denote the prevalence (%) of malignant masses in each centre. **Oncology centres** were: University Hospitals Leuven, Belgium (LBE); Università Cattolica del Sacro Cuore, Rome, Italy (RIT); Ospedale San Gerardo, Monza, Italy (OIT); General Faculty Hospital, Prague, Czech Republic (PCR); Istituto Europeo di Oncologia, Milan, Italy (CIT); Medical University Lublin, Poland (LPO); Karolinska University Hospital, Stockholm, Sweden (SSW); Skåne University Hospital Lund, Sweden (LSW); Università degli Studi di Udine, Italy (UDI); Istituto Nazionale dei Tumori, Naples, Italy (GIT); Università di Bologna, Italy (BIT). **Non-oncology centres** were: Skåne University Hospital Malmö, Sweden (MSW), Ziekenhuis Oost-Limburg, Genk, Belgium (GBE); Ospedale San Giovanni di Dio, Cagliari, Italy (SIT); DSC L. Sacco, University of Milan, Italy (MIT); Università degli Studi di Napoli, Naples, Italy (NIT); Institut Universitari Dexeus, Barcelona, Spain (BSP); Ospedale dei Bambini Vittore Buzzi, Milan, Italy (FIT).

Figure S2. The area under the receiver operating characteristic curve (AUC) with 95% CI for logistic regression model 2 (LR2) per contributing centre and for all centres combined using a meta-analysis approach and pooled data.



NC, not computable; CI, confidence interval. Numbers in brackets denote the prevalence (%) of malignant masses in each centre. **Oncology centres** were: University Hospitals Leuven, Belgium (LBE); Università Cattolica del Sacro Cuore, Rome, Italy (RIT); Ospedale San Gerardo, Monza, Italy (OIT); General Faculty Hospital, Prague, Czech Republic (PCR); Istituto Europeo di Oncologia, Milan, Italy (CIT); Medical University Lublin, Poland (LPO); Karolinska University Hospital, Stockholm, Sweden (SSW); Skåne University Hospital Lund, Sweden (LSW); Università degli Studi di Udine, Italy (UDI); Istituto Nazionale dei Tumori, Naples, Italy (GIT); University of Bologna, Italy (BIT). **Non-oncology centres** were: Skåne University Hospital Malmö, Sweden (MSW), Ziekenhuis Oost-Limburg, Genk, Belgium (GBE); Ospedale San Giovanni di Dio, Cagliari, Italy (SIT); DSC L. Sacco, University of Milan, Italy (MIT); Università degli Studi di Napoli, Naples, Italy (NIT); Institut Universitari Dexeus, Barcelona, Spain (BSP); Ospedale dei Bambini Vittore Buzzi, Milan, Italy (FIT).

Figure S3. The area under the receiver operating characteristic curve (AUC) with 95% CI for the Risk of Malignancy Index (RMI) per contributing centre and for all centres combined using a meta-analysis approach and pooled data.



NC, not computable; CI, confidence interval. Numbers in brackets denote the prevalence (%) of malignant masses in each centre. **Oncology centres** were: University Hospitals Leuven, Belgium (LBE); Università Cattolica del Sacro Cuore, Rome, Italy (RIT); Ospedale San Gerardo, Monza, Italy (OIT); General Faculty Hospital, Prague, Czech Republic (PCR); Istituto Europeo di Oncologia, Milan, Italy (CIT); Medical University Lublin, Poland (LPO); Karolinska University Hospital, Stockholm, Sweden (SSW); Skåne University Hospital Lund, Sweden (LSW); Università degli Studi di Udine, Italy (UDI); Istituto Nazionale dei Tumori, Naples, Italy (GIT); University of Bologna, Italy (BIT). **Non-oncology centres** were: Skåne University Hospital Malmö, Sweden (MSW), Ziekenhuis Oost-Limburg, Genk, Belgium (GBE); Ospedale San Giovanni di Dio, Cagliari, Italy (SIT); DSC L. Sacco, University of Milan, Italy (MIT); Università degli Studi di Napoli, Naples, Italy (NIT); Institut Universitari Dexeus, Barcelona, Spain (BSP); Ospedale dei Bambini Vittore Buzzi, Milan, Italy (FIT).

Sensitivity analysis for missingness of serum CA125

(AUC, area under the receiver operating characteristic curve; sens, sensitivity; spec, specificity; CI, confidence interval; LR2, logistic regression model2; RMI, risk of malignancy index; SR, simple rules; SR+mal, simple rules as first stage test with all cases in which the simple rules cannot be applied classified as malignant; SR+SA, simple rules as first stage test and all cases in which the simple rules cannot be applied classified by subjective assessment)

Table S5. Diagnostic test performance of LR2, Risk of Malignancy Index (RMI) and Simple Rules strategies using only complete cases for CA125 (n=1451).

	AUC (95% CI)	Sens (95% CI)	Spec (95% CI)
LR2	0.907 (0.886-0.925)	91.8 (88.8-94.1)	73.2 (63.6-81.0)
RMI	0.856 (0.820-0.886)	72.6 (68.8-76.2)	86.3 (80.9-90.3)
SR+Mal	/	96.7 (93.9-98.3)	67.1 (58.6-74.6)
SR+SA	/	93.6 (90.4-95.7)	84.6 (77.7-89.7)

Table S6. Diagnostic test performance of LR2, Risk of Malignancy Index (RMI) and Simple Rules strategies using only data from centres (N=9) with almost complete information (0-6% missingness) on CA125 (n=994).*

	AUC (95% CI)	Sens (95% CI)	Spec (95% CI)
LR2	0.924 (0.899-0.943)	91.6 (87.5-94.5)	76.7 (67.3-84.0)
RMI	0.867 (0.826-0.900)	72.1 (65.5-77.8)	88.1 (81.0-92.8)

SR+Mal	/	98.1 (86.3-99.8)	70.4 (58.2-80.2)
SR+SA	/	92.8 (86.3-96.4)	88.2 (82.9-92.1)

* A single imputation technique was used to impute missing values for CA125.

Table S7. Histology in cases with and without information on CA125

Histology	CA125 available n (%) ^a	CA125 missing, n (%) ^a	Total (n)
Endometrioma	161 (11.10)	183 (19.22)	344
Teratoma	109 (7.51)	122 (12.82)	231
Simple cyst + parasalpingeal cyst	56 (3.86)	50 (5.25)	106
Functional cyst	16 (1.10)	24 (2.52)	40
Hydrosalpinx + salpingitis	19 (1.31)	28 (2.94)	47
Peritoneal pseudocyst	5 (0.34)	13 (1.37)	18
Abscess	9 (0.62)	8 (0.84)	17
Fibroma	79 (5.44)	51 (5.36)	130
Serous cystadenoma	156 (10.75)	103 (10.82)	259

Histology	CA125 available n (%) ^a	CA125 missing, n (%) ^a	Total (n)
Mucinous cystadenoma	99 (6.82)	84 (8.82)	183
Rare benign	31 (2.14)	17 (1.79)	48
Primary invasive stage I	97 (6.69)	31 (3.26)	128
Primary invasive stage II	35 (2.41)	12 (1.26)	47
Primary invasive stage III	306 (21.09)	91 (9.56)	397
Primary invasive stage IV	46 (3.17)	15 (1.58)	61
Rare primary invasive	41 (2.83)	27 (2.84)	68
Borderline stage I	91 (6.27)	44 (4.62)	135
Borderline stage II	6 (0.41)	0 (0.00)	6
Borderline stage III	11 (0.76)	1 (0.11)	12
Metastatic	78 (5.38)	48 (5.04)	126
Total	1451	952	2403

^aPercentages have been calculated per column