

1 **TITLE PAGE:**

2 Full title of the paper: Reproducibility of the Endometriosis Fertility Index: a prospective
3 inter/intra-rater agreement study

4 Authors: C. Tomassetti^{1, 2*}, C. Bafort¹, C. Meuleman^{1, 2}, M. Welkenhuysen¹, S. Fieuws³, T.
5 D'Hooghe²

6 ¹ Department of Obstetrics and Gynaecology, Leuven University Fertility Centre, University
7 Hospitals Leuven, Herestraat 49, Leuven 3000, Belgium

8 ² Department of Development and Regeneration, KU Leuven, Herestraat 49, Leuven 3000,
9 Belgium

10 ³ Department of Public Health, Interuniversity Centre for Biostatistics and Statistical
11 Bioinformatics, KU Leuven, Kapucijnenvoer 35, Leuven 3000, Belgium.

12 * Correspondence address. E-mail: carla.tomassetti@uzleuven.be

13 Shortened running title: Reproducibility of the EFI

14 **ABSTRACT:**

15 Objective: To evaluate the reproducibility of the EFI (Endometriosis Fertility Index).

16 Design: Single-cohort prospective observational study.

17 Setting: University hospital.

18 Population: Women undergoing laparoscopic resection of any rASRM-stage endometriosis.

19 Methods: Details of pre- and per-operative findings were collected into a coded research
20 file. EFI-scoring was performed 'en-bloc' by three different raters (expert-1 (C.T.), expert-2
21 (C.M.), junior (C.B.)). Required sample size: 71. Definitions used for agreement: clinical
22 (scores within same range: 0-4, 5-6, 7-10) and numerical (difference ≤ 1 EFI-point).

23 Main outcome measures: Primary outcome: rate of clinical agreement between two experts.
24 Secondary outcomes: expert numerical agreement, clinical and numerical agreement
25 between expert-1 and junior and within expert-1 (intra-observer), agreement of rASRM-
26 score and -stage.

27 Results: A near-to-perfect 'inter-expert' clinical agreement rate (1.000 (95% CI 0.956-1.000),
28 $p=0.0149$) was observed. The numerical agreement between two experts was also high
29 (0.988 (95% CI 0.934-1.000)); similarly high agreement rates were observed for both 'junior-
30 expert' comparison (clinical 0.963 (95% CI 0.897-0.992), numerical 0.988 (95% CI 0.934-
31 1.000) and 'intra-expert' comparisons (clinical 0.988 (95% CI 0.934-1.000); numerical 1.000
32 (95% CI 0.956-1.000)). Reasons for disagreements were different scoring of the least-
33 function score and disagreements in rASRM-scores. The reproducibility of the rASRM-score
34 was clearly inferior to that of the EFI for all comparisons.

35 Conclusion: The EFI can be reproduced reliably by different raters, further supporting its use
36 in daily clinical practice as the principal clinical tool for postoperative fertility
37 counselling/management of women with endometriosis.

38 Funding: FWO (Fund for Scientific Research, Flanders)

39 Keywords: Endometriosis, laparoscopy, infertility, classification, reproducibility of results

40 **TWEETABLE ABSTRACT:**

41 A study confirming the high reproducibility of the EFI substantiates its use in daily clinical
42 practice.

43 INTRODUCTION

44 Although the rASRM (revised American Society for Reproductive Medicine) score¹ is the
45 most frequently used surgical staging system for endometriosis to date, it has some serious
46 limitations. First, its reproducibility has only been described as being 'fair to good'²⁻⁵, thus
47 prone to inter-observer variability. Second, it is not effective for predicting clinical outcomes
48 of treatment, especially pregnancy rates in infertile patients.⁶⁻⁸ For the latter reason, in 2010
49 Adamson and Pasta developed the EFI (Endometriosis Fertility Index), which now is a
50 thoroughly validated scoring system that predicts pregnancy rates without using ART
51 (assisted reproductive technology) treatment in postoperative endometriosis patients who
52 suffer from infertility and takes into account all endometriosis rASRM stages.⁹⁻¹³
53 Consequently, the EFI has been adopted by the WES (World Endometriosis Society) in their
54 consensus on the classification of endometriosis.¹⁴ In the EFI, 5 out of 10 possible points are
55 based on patient characteristics such as age, duration of infertility and history of pregnancy.
56 Parts of the rASRM staging account for 2 points of the EFI. Being an end-of-surgery staging,
57 the rest of the score is based on visual observation and qualitative assessment by the
58 surgeon (adnexal 'least function' score: 3 points). Especially the surgical part of the EFI
59 score could make it prone to differences in interpretations by different observers, which in
60 turn could have an effect on subsequent patient management. In the paper by Adamson and
61 Pasta⁹ who developed the EFI, a sensitivity analysis was reported to assess the effect on the
62 EFI of potentially assumed differences in the assignment of the adnexal least function score
63 by different surgeons, it was concluded that an EFI change of more than 1 point would only
64 be present in 5.4% of the cases; the authors further stated that changes in the EFI would be
65 material only for the middle values. However, this was only a theoretical exercise, and a
66 possible added influence of the poor inter-observer agreement of the rASRM score and

67 stage was not accounted for. Also, to our knowledge, no true inter-observer
68 variability/reliability assessment for the EFI has been performed so far.
69 The objective of this study was to evaluate whether the EFI score can be reproduced reliably
70 by different raters, i.e. whether the inter-observer variability is absent or low enough to
71 avoid a relevant impact on clinical patient management. Additionally, intra-observer
72 agreement of the EFI, and inter- and intra-observer agreement on the rASRM score were
73 also studied.

74

75 **METHODS**

76 **Study design**

77 This is a single cohort prospective observational (non-interventional) study in **women**
78 scheduled for endometriosis surgery of any rASRM stage at the LUFC (Leuven University
79 Fertility Centre) of the University Hospitals Leuven (**UZ Leuven**, Belgium). The study was
80 conducted, based on patient data gathered from surgical procedures performed from June
81 13th, 2016 until December 22nd, 2016 included. Three assessors with a different profile were
82 chosen: C.T. is an expert surgeon with a long experience of EFI-scoring, C.M. is also an expert
83 surgeon who only occasionally uses the EFI score, and C.B. is a trainee in obstetrics and
84 gynaecology.

85 Three different comparison levels were decided when designing the study protocol:
86 comparison between rating of expert 1 (C.T.) and expert 2 (C.M.) (further referred to as
87 'inter-expert'), between rating of expert 1 (C.T.) and junior (C.B.) ('junior-expert'), and
88 between rating of the first and the second session of expert 1 (C.T.) ('intra-expert').

89 The choice of experts as well as a trainee makes this study interesting not only for a tertiary
90 referral centre for endometriosis, but also for those with less experience with the disease

91 (such as trainees). There was no involvement from patients or public in the development of
92 this study.

93

94 **Study population – eligibility criteria**

95 The LUFC is a tertiary referral centre for both endometriosis and reproductive medicine.

96 Women of the reproductive age group (18-45 years), undergoing CO₂-laser laparoscopic
97 surgery at the LUFC for diagnosis and treatment of endometriosis, with confirmed diagnosis
98 on pathological examination, were eligible for this study. Indication for surgery had to be at
99 least one of the following: infertility of \geq 12 months, clinical examination and/or pain
100 symptoms suggesting endometriosis, ultrasound (and/or other relevant imaging) findings
101 suggesting endometriosis, previous surgical diagnosis of endometriosis. Laparoscopic
102 procedures in the setting of a day surgery centre as well as a hospitalization setting were
103 included. Patients were excluded in case they had a history of or were planned to undergo a
104 hysterectomy and/or bilateral salpingo-oophorectomy, if endometriosis lesions were not
105 completely resected (e.g. only marsupialization of an endometrioma), if photographic
106 documentation was not performed or not compatible with study quality standards (see
107 description in study procedures), or if informed consent was not obtained.

108

109 No extra study-related patient informed consent was necessary, since patients agreed
110 preoperatively in their surgical informed consent form that their clinical data (which
111 routinely include photographic documentation of the surgery) may be stored and used for
112 scientific purposes. Confidentiality was ascertained by anonymously transferring the
113 necessary patient data into a specifically designed research file (CRF or case report form).

114

115 **Data recording and procedures**

116 Next to demographical and clinical data (including results from clinical examination, imaging,
117 extensive surgical reports and those specific data necessary for calculation of the historical
118 part of the EFI), standardised photographic documentation of the laparoscopic findings was
119 done, both at the start and **at the** end of the surgery as per WERF-EPHect-guidelines.¹⁵

120 Although no video recordings were used, the mobility of the tube and ovary was be assessed
121 on photograph by lifting the adnexa out of the ovarian fossa.

122

123 All necessary data were transferred to the CRF by C.B., a second-year obstetrics and
124 gynaecology resident-in-training at the time of the study. In this CRF, data were anonymized
125 and standardized, information on date of surgery was removed, and a unique and
126 anonymous study number was allocated to each patient, to guarantee confidentiality and
127 blinding of the assessors.

128 Surgical procedures were performed by C.T. or C.M., **who are** both reproductive
129 endocrinologists as well as reproductive gynaecological surgeons with a specific expertise in
130 the treatment of all forms of **extensive**-endometriosis.¹⁶

131

132 Only when the appropriate sample size was reached and subsequently all CRFs had been
133 created, 'en-bloc' rating sessions were organized for each rater. All raters scored the EFI
134 based on all the information in the CRF separately and independently from each other.

135 Completed scoring forms were kept under lock by the study coordinator until the time of
136 data analysis. There was at least four weeks between the last surgical procedure and the
137 first rating session. Recall or other bias of the raters was avoided due to the time interval
138 between surgery and rating session, the anonymization of the patient information in the

139 CRF, the different order in which the files were rated, and the closed storage of the
140 completed scoring forms.

141 During the rating session, all raters completed two scoring forms per patient: one for the
142 rASRM and one for the surgical part of the EFI, based on the pre- and per-operative
143 information in the CRF. Four weeks after her first rating, C.T. repeated the rating session for
144 intra-observer variability assessment. Since the historical EFI factors are not prone to be
145 interpreted differently by different observers, they were filled directly into the final study
146 database but weren't scored by each rater separately. For the final calculation of the total
147 EFI score for each patient and for each rater/session, the (fixed) historical and (differentially
148 rated) surgical EFI points were added together in the study database.

149

150 **Outcomes**

151 The primary outcome studied was the percentage of clinical agreement of the EFI-score in
152 the 'inter-expert' comparison. Clinical agreement was defined as having no impact on the
153 subsequent clinical decision pathway regarding fertility management as currently used at
154 the LUFC, meaning that EFI-scores should be within the same range (low EFI range: 0-4,
155 median EFI range: 5-6, high EFI range: 7-10).

156 Secondary outcomes studied were: clinical agreement on the EFI-score for 'junior- expert'
157 and 'intra-expert' comparison, numerical agreement on the EFI-score (defined as a
158 maximally allowed absolute difference in EFI-score of 1 point, regardless of the above
159 mentioned range) and agreement on rASRM-score/stage for all three comparisons ('inter-
160 expert', 'junior-expert', 'intra-expert').

161

162 **Sample size estimation**

163 This study was designed to show that the percentage of agreement between two senior
164 raters (inter-expert comparison) is higher than 95% for clinical agreement (primary
165 outcome). Based on a one-sided binomial test for a single proportion with $\alpha=0.05$,
166 expecting the true percentage of discrepancies to be $<0.001\%$, the minimal sample size
167 equals 71 subjects to have at least 80% power to show that the percentage of discrepancies
168 is lower than 5%. The minimally required sample size was therefore set at 71.

169

170 **Statistical analysis**

171 A one-sided binomial test with $\alpha=0.05$ for a single proportion was used to test if the
172 observed proportion of clinical agreement between both experts was significantly higher
173 than 95%. For all percentages of agreement, two-sided 95% CIs are reported as well.

174 Weighted kappas (with the classical quadratic weighing), which are widely used in
175 agreement studies¹⁷⁻²¹, were reported both for the total EFI and for the rASRM stage, where
176 a kappa of 1 indicates perfect agreement and 0 indicates agreement equivalent to chance.

177 Bland-Altman plots were used to visualize the agreement of the total rASRM score.²² Such
178 plots provide information on the bias (the mean difference as tested with a paired t-test),
179 the expected range of the difference in scores (95% LOA (limits of agreement)) and the
180 possible dependency of the difference on the level of the score. Additionally, the ICC (intra-
181 class correlation coefficient) was given for the quantification of the agreement for the total
182 rASRM score.²³

183 All analyses have been performed using SAS software, version 9.4 of the SAS System for
184 Windows.

185

186 **RESULTS**

187 156 patients underwent laparoscopic surgery at the LUFC between June 13th, 2016 until
188 December 22nd, 2016 included. 29 patients did not have endometriosis at laparoscopy. Out
189 of 127 laparoscopies for endometriosis, 10 did not fit the inclusion criteria: 2 patients were
190 outside age range, 3 had incomplete surgery for the pelvis, 4 underwent planned 2-step
191 surgery and 1 patient had **additional** other pathology ~~in addition to endometriosis~~. Out of
192 the 117 eligible patients, 35 did not have sufficiently detailed photographic documentation,
193 so finally 82 patients were included for creation of CRFs, rating and analysis, which was more
194 than the minimally required sample size. Among the included patients, 41 surgical
195 procedures were performed by C.T., and 41 by C.M.; 13 were assisted by C.B..

196

197 ~~Prior to surgery, the most frequently reported endometriosis-related pain symptom was~~
198 ~~dysmenorrhea (75/82, 91.5%), including mostly moderate (25/75; 33.3%) or severe (41/75;~~
199 ~~54.7%) dysmenorrhea. Other prevalent baseline symptoms included dyschezia (45/82;~~
200 ~~54.9%), and/or rectal bleeding (16/82; 19.5%), deep dyspareunia (37/81; 45.7%), chronic~~
201 ~~pelvic pain (36/82; 43.9%), and mictalgia (24/82; 29.3%). In addition, 39/82 (47.5%) had a~~
202 ~~history of diagnostic and/or incomplete therapeutic surgery for endometriosis, and 15/82~~
203 ~~(18.29%) and 13/82 (15.85%) patients had a history of treatment with IUI or ART~~
204 ~~respectively.~~

205

206 **EFI**

207 Baseline demographic characteristics, including those necessary for calculation of the
208 historical points of the EFI ~~and the type of endometriosis lesions identified~~, are shown in
209 Table 1. The most frequently found type of endometriosis lesions were peritoneal implants

210 (78/82, 90.2%), followed by deep (64/82, 78.1%), superficial ovarian 42/82 (51.2%) and
211 cystic ovarian (23/82, 28%) endometriosis.

212 Table 2A shows the results for EFI score agreement according to both definitions described
213 above, and the weighted kappa for the 3 comparisons made. The majority of included
214 patients had high scores for the historical part of the EFI (4 points, 45/82 (54.88%) or 5
215 points 23/82 (28.05%)), as reflected partly in the clustering of the higher EFI-scores (Table 3).
216 This is comparable with a previous study in our population ¹⁰, which confirms the studied
217 population as representative for our clinic.

218

219 ***Inter-expert EFI comparison***

220 For the 'inter-expert' clinical agreement, the study hypothesis was confirmed, namely that
221 the rate of agreement was higher than 95%, which was near-to-perfect (1.000 (95% CI 0.956-
222 1.000), one-sided p-value=0.0149).

223 The 'inter-expert' numerical agreement was slightly lower than the clinical agreement (with
224 the lower limit of the 95% CI just below 0.95: 0.988 (95% CI 0.934-1.000)).

225

226 Table 3 shows the details of agreement for the 'inter-expert' comparison (similar data on the
227 other comparisons can be supplied upon request). In 9 cases, EFI scores did not reach
228 absolute agreement between both experts C.T. and C.M., of which only 1 led to the defined
229 'numerical disagreement' (EFI score 4 versus 2). Out of these 9 cases, 3 were due to
230 differences in rASRM score (1 in lesion score ≤ 16 , 2 in total score ≤ 71), and 6 were
231 due to C.T. giving a lower LF score than C.M. (4 with bilateral vaporization of superficial
232 ovarian endometriosis, 1 with treatment of an endometrioma, and 1 for of tubal/fimbrial
233 functionality).

234

235 Junior-Expert EFI comparison

236 For the comparison 'junior-expert', in general the rate of agreement was slightly lower than
237 for the inter-expert EFI comparison, but still around 90% or more when taking into account
238 the lower limit of the 95% CI (0.963 (95% CI 0.897-0.992) for clinical agreement, 0.988 (95%
239 CI 0.934-1.000) for numerical agreement).

240 Details of disagreement were as follows: 1 case with both numerical and clinical
241 disagreement and 2 cases with clinical disagreement only, out of the total of 15/82 files with
242 any difference in EFI scoring between junior and expert. Of these 15 cases, 4 were due to a
243 difference in total rASRM score (> or ≤71), 7 due to different ovarian LF score (of which 1 led
244 to clinical disagreement) and 4 due to different tubal/fimbrial LF score (of which 1 led to
245 clinical, and 1 to clinical and numerical disagreement).

246

247 Intra-expert EFI comparison

248 Agreement was also high for the 'intra-expert' comparison (numerical agreement (1.000
249 (95% CI 0.956-1.000), clinical agreement (0.988 (95% CI 0.934-1.000))).

250 For this comparison, only 1 case had clinical disagreement out of a total of 7/82 of cases with
251 any difference in EFI score. Of these latter 7 cases, 1 difference was attributed to the total
252 rASRM score, and 6 to the LF score (4 on ovarian function and 2 on tubal/fimbrial function
253 (amongst which 1 led to clinical disagreement)).

254

255 rASRM scoring and staging

256 From Figure 1, showing the Bland-Altman plot and statistical analysis of the agreement on
257 the total rASRM score (in points), it's clear that the variability for the total rASRM score

258 given is very large for all 3 comparisons. Indeed, although the mean differences of assigned
259 rASRM points may be small (confirming a low risk for fixed bias), their SDs are large, and the
260 95% LOA (limits of agreement) span a width of 40 points or more, which is comparable to 4
261 rASRM stages.

262 Table 2B describes the analysis of agreement on rASRM stage, explained by rate of
263 agreement and weighted kappa; these results are consistently lower than those obtained for
264 the EFI (Table 2A). The **supplemental figure S1** shows an example of a **woman** where
265 complete agreement between all raters was found.

266

267 **Relationship between rASRM and EFI**

268 **Supplemental Figure S2 shows a boxplot of the distribution of rASRM total score for each EFI**
269 **range (for expert 1). This illustrates that in general there was a negative correlation between**
270 **the rASRM (points/stage) and EFI range. Interestingly, 43/62 (72,58%) of women with a high**
271 **EFI also have rASRM stage III-IV endometriosis (Figure S1B).**

272

273 **DISCUSSION**

274 Our study represents the first report on inter- and intra-observer reproducibility of the EFI
275 and demonstrates high intra- and inter-agreement rate with narrow 95% CIs. More
276 specifically, we **have** confirmed our hypothesis that clinical agreement for the 'inter-expert'
277 comparison (primary outcome) was higher than 95%. These results concur with the
278 hypothetical assumption based on the sensitivity analysis on the EFI by Adamson and Pasta⁹,
279 **as explained in the introduction**. In addition, very high agreements were also reported for
280 numerical 'inter-expert' agreement, clinical and numerical 'junior-expert' and 'intra-expert'
281 comparison (secondary outcomes), **we found very high rates** although not near-to-perfect as

282 for clinical “inter-expert” agreement. In other words, the high reproducibility supports the
283 use of the EFI in daily clinical practice as a very relevant clinical tool for management and
284 counselling of postoperative endometriosis patients on their reproductive outcome.

285

286 Disagreement between raters could be largely explained by differential rating of the least
287 function score, and of the rASRM score. The influence of lower reproducibility of the rASRM
288 score on the EFI score reproducibility was not taken into account in the sensitivity analysis
289 by Adamson and Pasta⁹ but is now identified in our data – next to the least function score –
290 as a potential weak spot in the reproducibility of the EFI score.

291

292 Our study was designed to avoid bias in several ways. First of all, the assessment of the EFI
293 was done based on a combination of patient history information, standardized operative
294 reports and complete photographic series of the operative site, in order to prevent any
295 misclassification of rASRM staging and associated adnexal adhesions as much as possible.^{5, 21}
296 Second, to blind raters to the personal details of patients, a coded CRF was used for rating
297 instead of the patient file itself. **Third, to avoid recall bias, a standardized and anonymized**
298 **CRF was used. Additionally, ‘en-bloc’ rating sessions, with random order of patient files,**
299 **were organised for each rater. Fourth,** since C.T. had the most experience in calculating the
300 EFI in clinical practice, her first rating was therefore chosen as standard to assess agreement
301 with the second expert (‘inter-expert’), the junior surgeon (‘junior-expert’) and within one
302 rater (‘intra-expert’).

303

304 Out of the 117 eligible patients, 35 were excluded because they did not have sufficiently
305 detailed photographic documentation. This was not considered as a flaw, but merely a

306 consequence of the fact that the study was conducted in a real life turbulent clinical setting
307 (different surgeons, different operation theatres, technical difficulties etc.). Patients files
308 were only included if photographic documentation (both pre- and postoperative) met the
309 criteria as defined per WERF-EpHect procedures.¹⁵ Despite this strict selection, our study
310 population was still representative for the population in our clinic (see result section), and
311 the minimally required sample size was more than met.

312

313 For all comparisons made, the rate of agreement was lower for the rASRM endometriosis
314 total score and rASRM endometriosis stage than for EFI score, despite our efforts to avoid
315 misclassification as described above. With respect to assessments of rASRM total score, the
316 width of variation was very high, and therefore the finding of a low mean error for all three
317 comparisons is not necessarily reassuring. Indeed, also ICCs are falsely inflated, since they
318 compare the difference within a subject to the difference between subjects, and in a more
319 uniform population (where the range of rASRM total score would be smaller than in our
320 population) the ICC would be considerably lower if still similar variation between observers
321 would be found.

322 With respect to rASRM stage assessment, agreements were also lower than for the EFI, as
323 explained by the lower values for weighted kappa and the lower limits of 95% CI for
324 agreement per se. When comparing results for weighted kappa, it should be noted that, in
325 contrast to the EFI where 11 possible categories are withheld (0-10, including both), in the
326 rASRM classification only 4 stages are categorized, but still results on rASRM stage showed a
327 markedly higher variability.

328

329 This study has a number of limitations that should be taken into account. First, the
330 relatively small numbers of raters involved may be a negative point, although this was
331 accounted for in the sample size calculation as discussed in the methodology section.
332 Second, raters with various levels of expertise of EFI scoring were included as describe in the
333 methodology section. The junior rater was also trained by the expert rater amongst others.
334 Therefore, we would suggest future studies on the reproducibility on the EFI to include a
335 larger number of observers and a more varied pool of observers preferably from various
336 centres with different expertise. Third, risk of recall bias cannot completely be excluded
337 since both experts performed all laparoscopies, and the junior assisted some procedures.
338 Fourth, the use of photographic documentation only rather than video recording during the
339 surgery to assess both the initial endometriosis lesions and the least function score at
340 conclusion of surgery may be less precise. However, as per WERF-EPHect-guidelines¹⁵,
341 photographic documentation only was assumed to be sufficient for the aim of our study and
342 could easily be embedded in our daily clinical practice. Fifth, next to photographic
343 documentation, standardized operative reports were provided to the raters, which could
344 positively influence the precision of the rating as described in the inter-rater agreement
345 study of Schliep et al²¹. However, this argument can easily be rejected since – in contrast to
346 the EFI – the reproducibility of rASRM score and stage remained poor. Finally, the estimated
347 sample size for the primary outcome (i.e. percentage of clinical agreement) may appear too
348 small, although the null hypothesis was derived from the EFI development study⁹. In
349 hindsight, the assumption used in the calculation (true percentage of discrepancies lower
350 than 0.001%) could be considered as too optimistic.

351

352 **CONCLUSION**

353 In addition to already vast evidence confirming the EFI score to be superior to the rASRM
354 score/stage for the prediction of reproductive outcome after surgery, our study has now
355 clearly demonstrated that EFI scoring is highly reproducible. This high reproducibility is far
356 better than for the rASRM scoring/staging, even for a trainee. Collectively, this evidence
357 supports the standard use of the EFI score next to the rASRM score/stage in daily clinical
358 practice as also advised by the WES ¹⁴, and the replacement of the rASRM stage/score by the
359 EFI score for postoperative fertility counselling of endometriosis patients. Preferably, our
360 data on reproducibility of the EFI score, as presented in this study, should be confirmed by
361 other groups, ideally by using a similar methodology but with a larger number of raters to
362 enhance comparability with our data.

363

364 **ACKNOWLEDGEMENTS**

365 None

366

367 **DISCLOSURE OF INTERESTS**

368 - The LUFC received unrestricted research grants from Ferring Pharmaceuticals and
369 Merck SA.

370 - C.T. received consultancy fees from Gedeon Richter and Lumenis (payed to UZ/KU
371 Leuven, no private revenue), sponsoring from Ferring Pharmaceuticals, Merck SA,
372 Gedeon Richter and Bayer to travel to and attend scientific meetings.

373 - C.B., M.W. and S.F. report no conflicts of interest

374 - C.M. received consultancy fees from Lumenis and Merck SA (payed to KULeuven, no
375 private revenue)

376 - T.D. has become vice-president and head of global medical affairs infertility for the
377 multinational pharmaceutical company Merck (Darmstadt, Germany) from October
378 1st, 2015. He continues on a part time basis his academic appointment as Professor
379 of Reproductive Medicine at the University of Leuven (KU Leuven) in Belgium and as
380 Adjunct Professor at the Department of Obstetrics and Gynaecology at Yale
381 University, New Haven, CO, USA.

382 - More details can be found on the ICMJE forms for disclosure of potential conflicts
383

384 **CONTRIBUTION TO AUTHORSHIP**

385 C.T.: study design, study performance (surgery, CRF design, rating), article writing and
386 editing

387 C.B.: study performance (surgery, CRF data-transfer, rating), article editing

388 C.M.: study performance (surgery, rating), article editing

389 M.W.: study performance (CRF data-transfer, file coding), article editing

390 S.F.: statistical analysis, article editing

391 T.D.: study design, article editing

392

393

394 **DETAILS OF ETHICS APPROVAL**

395 This study was approved by the ethics committee of the UZ Leuven on June 8th, 2016

396 (internal UZ Leuven Trial Registration Number: S59221); competent authority approval was

397 not necessary since the study was observational.

398

399 **FUNDING**

400 This study was supported by FWO (Fund for Scientific Research, Flanders).

401 **REFERENCES**

- 402 1. American Society for Reproductive M. Revised American Society for Reproductive
403 Medicine classification of endometriosis: 1996. *Fertility and Sterility*. 1997;67(5):817-21.
- 404 2. Rock JA, Zoladex Endometriosis Study Group JA. The revised American Fertility Society
405 classification of endometriosis: reproducibility of scoring. *Fertility and Sterility*.
406 1995;63(5):1108-10.
- 407 3. Hornstein MD, Gleason RE, Orav J, Haas ST, Friedman AJ, Rein MS, et al. The
408 reproducibility of the revised American Fertility Society classification of endometriosis. *Fertil*
409 *Steril*. 1993;59(5):1015-21.
- 410 4. Lin SY, Lee RKK, Hwu YM, Lin MH. Reproducibility of the revised American Fertility Society
411 classification of endometriosis using laparoscopy or laparotomy. *International Journal of*
412 *Gynecology & Obstetrics*. 1998;60(3):265-9.
- 413 5. Schliep CK, Stanford BJ, Chen KZ, Zhang OB, Dorais WJ, Boiman Johnstone BE, et al.
414 Interrater and Intrarater Reliability in the Diagnosis and Staging of Endometriosis. *Obstetrics*
415 *& Gynecology*. 2012;120(1):104-12.
- 416 6. Palmisano GP, Adamson GD, Lamb EJ. Can staging systems for endometriosis based on
417 anatomic location and lesion type predict pregnancy rates? *Int J Fertil Menopausal Stud*.
418 1993;38(4):241-9.
- 419 7. Vercellini P, Fedele L, Aimi G, De Giorgi O, Consonni D, Crosignani PG. Reproductive
420 performance, pain recurrence and disease relapse after conservative surgical treatment for
421 endometriosis: the predictive value of the current classification system. *Human*
422 *Reproduction*. 2006;21(10):2679-85.
- 423 8. Adamson DG. Endometriosis classification: an update. *Current Opinion in Obstetrics and*
424 *Gynecology*. 2011;23(4):213-20.

- 425 9. Adamson GD, Pasta DJ. Endometriosis fertility index: the new, validated endometriosis
426 staging system. *Fertility and Sterility*. 2010;94(5):1609-15.
- 427 10. Tomassetti C, Geysenbergh B, Meuleman C, Timmerman D, Fieuws S, D'Hooghe T.
428 External validation of the endometriosis fertility index (EFI) staging system for predicting
429 non-ART pregnancy after endometriosis surgery. *Human Reproduction*. 2013;28(5):1280-8.
- 430 11. Zeng C, Xu J-N, Zhou Y, Zhou Y-F, Zhu S-N, Xue Q. Reproductive Performance after
431 Surgery for Endometriosis: Predictive Value of the Revised American Fertility Society
432 Classification and the Endometriosis Fertility Index. *Gynecologic and Obstetric Investigation*.
433 2014;77(3):180-5.
- 434 12. Garavaglia E, Pagliardini L, Tandoi I, Sigismondi C, Viganò P, Ferrari S, et al. External
435 Validation of the Endometriosis Fertility Index (EFI) for Predicting Spontaneous Pregnancy
436 after Surgery: Further Considerations on Its Validity. *Gynecologic and Obstetric Investigation*.
437 2015;79(2):113-8.
- 438 13. Boujenah J, Bonneau C, Hugues JN, Sifer C, Poncelet C. External validation of the
439 Endometriosis Fertility Index in a French population. *Fertility and Sterility*. 2015;104(1):119-
440 23.e1.
- 441 14. Johnson NP, Hummelshoj L, Adamson GD, Keckstein J, Taylor HS, Abrao MS, et al. World
442 Endometriosis Society consensus on the classification of endometriosis. *Human*
443 *Reproduction*. 2017;32(2):315-24.
- 444 15. Becker CM, Laufer MR, Stratton P, Hummelshoj L, Missmer SA, Zondervan KT, et al.
445 World Endometriosis Research Foundation Endometriosis Phenome and Biobanking
446 Harmonisation Project: I. Surgical phenotype data collection in endometriosis research.
447 *Fertility and Sterility*. 2014;102(5):1213-22.

- 448 16. Meuleman C, Tomassetti C, Wolthuis A, Van Cleynenbreugel B, Laenen A, Penninckx F, et
449 al. Clinical outcome after radical excision of moderate-severe endometriosis with or without
450 bowel resection and reanastomosis: A prospective cohort study. *Annals of Surgery*.
451 2014;259(3):522-31.
- 452 17. Viera AJ, Garrett JM. Understanding interobserver agreement: the kappa statistic. *Fam*
453 *Med*. 2005;37(5):360-3.
- 454 18. Holland TK, Hoo WL, Mavrellos D, Saridogan E, Cutner A, Jurkovic D. Reproducibility of
455 assessment of severity of pelvic endometriosis using transvaginal ultrasound. *Ultrasound*
456 *Obstet Gynecol*. 2013;41:210-5.
- 457 19. Zahn CM, Luigi KFR, Olsen C, Whitworth SA, Washington A, Crothers B. Reproducibility
458 of Endocervical Curettage Diagnosis. *Obstet Gynecol*. 2011;118:240-8.
- 459 20. Paternot G, Devroe J, Debrock S, D'Hooghe TM, Spiessens C. Intra- and inter-observer
460 analysis in the morphological assessment of early-stage embryos. *Reproductive Biology and*
461 *Endocrinology*. 2009;7:105.
- 462 21. Schliep K, Chen Z, Stanford J, Xie Y, Mumford S, Hammoud A, et al. Endometriosis
463 diagnosis and staging by operating surgeon and expert review using multiple diagnostic
464 tools: an inter-rater agreement study. *BJOG: An International Journal of Obstetrics &*
465 *Gynaecology*. 2017;124(2):220-9.
- 466 22. Bland JM, Altman DG. Measuring agreement in method comparison studies. *Statistical*
467 *Methods in Medical Research*. 1999;8(2):135-60.
- 468 23. McGraw KO, Wong SP. Forming Inferences About Some Intraclass Correlation
469 Coefficients. *Psychological Methods*. 1996;1(1):30-46.

470 **Table 1:** Baseline characteristics, including historical factors of the EFI and their translation
 471 into EFI-points, for the total population (N=82) (NA = not applicable)

Characteristic	Mean ± SD	Median (IQR)	Number of patients/total (%)
Pain symptoms	NA	NA	
- Dysmenorrhea			75/82 (91,5%)
- Dyschezia			45/82 (54,9%)
- Rectal bleeding			16/82 (19,5%)
- Deep dyspareunia			37/81 (45,7%)
- Chronic pelvic pain			36/82 (43,9%)
- Mictalgia			24/82 (29,3%)
History of diagnostic/incomplete surgery	NA	NA	39/82 (47,5%)
History of fertility treatment	NA	NA	
- IUI			15/82 (18,29%)
- ART			13/82 (15,85%)
Age (in years)	31.5 ± 4.65	31.2 (28.4-34.8)	0 EFI points (age 40+): 1/82 (1.22%) 1 EFI point (age 36-39): 16/82 (19.51%) 2 EFI points (age <36): 65/82 (79.27%)
Duration of infertility (in months)	17.1 ± 22.17	13.0 (0-29)	0 EFI points (>3 years): 7/82 (8.54%) 1 EFI point (≤3 years): 75/82 (91.46%)
Prior pregnancy	NA	NA	0 EFI point (never): 49/82 (59.76%) 1 EFI point (ever): 33/82 (40.24%)
EFI: total historical points	NA	NA	0 EFI points: 0/82 (0%) 1 EFI point: 1/82 (1.2%) 2 EFI points: 6/82 (7.3%) 3 EFI points: 7/82 (8.5%) 4 EFI points: 45/82 (54.5%) 5 EFI points: 23/82 (28.1%)

472

473

474 **Table 2:** Agreement for total EFI score and rASRM stage between raters

475 Table 2A: Agreement for total EFI score between raters

Comparison	Clinical agreement EFI score (95% CI)	Numerical agreement EFI score (95% CI)	Weighted kappa EFI score (95% CI)
Inter-expert	1.000 (0.956-1.000) *	0.988 (0.934-1.000)	0.942 (0.904-0.980)
Junior-expert	0.963 (0.897-0.992)	0.988 (0.934-1.000)	0.907 (0.858-0.956)
Intra-expert	0.988 (0.934-1.000)	1.000 (0.956-1.000)	0.959 (0.929-0.990)

476 *primary outcome: one-sided p-value = 0.0149

477 Table 2B: Analysis of (absolute) agreement on rASRM stage for the different comparisons

Comparison	Agreement rASRM stage (95% CI)	Weighted kappa rASRM stage (95% CI)
Inter-expert	0.841 (0.744-0.913)	0.752 (0.621-0.882)
Junior-expert	0.890 (0.802-0.949)	0.752 (0.721-0.882)
Intra-expert	0.915 (0.832-0.965)	0.907 (0.847-0.968)

478

479

480 **Table 3:** Cross-tabulation of the frequency of a given EFI-score for the inter-expert
 481 comparison – raw data (note that no score below 2 was given by any of the two raters).

		EFI by expert 2									Total
		2	3	4	5	6	7	8	9	10	
EFI by expert 1	2	1	0	1	0	0	0	0	0	0	2
	3	0	2	0	0	0	0	0	0	0	2
	4	0	0	3	0	0	0	0	0	0	3
	5	0	0	0	7	0	0	0	0	0	7
	6	0	0	0	0	6	0	0	0	0	6
	7	0	0	0	0	0	12	2	0	0	14
	8	0	0	0	0	0	2	13	4	0	19
	9	0	0	0	0	0	0	0	19	0	19
	10	0	0	0	0	0	0	0	0	10	10
	Total	1	2	4	7	6	14	15	23	10	82

482