

This is an Open Access document downloaded from ORCA, Cardiff University's institutional repository: <https://orca.cardiff.ac.uk/id/eprint/114345/>

This is the author's version of a work that was submitted to / accepted for publication.

Citation for final published version:

Leijdekkers, J. A., Eijkemans, M. J. C., van Tilborg, T. C., Oudshoorn, S. C., McLernon, D. J., Bhattacharya, Siladitya, Mol, B. W. J., Broekmans, F. J. M. and Torrance, H. L. 2018. Predicting the cumulative chance of live birth over multiple complete cycles of in vitro fertilization: an external validation study. *Human Reproduction* 33 (9), pp. 1684-1695. [10.1093/humrep/dey263](https://doi.org/10.1093/humrep/dey263)

Publishers page: <http://dx.doi.org/10.1093/humrep/dey263>

Please note:

Changes made as a result of publishing processes such as copy-editing, formatting and page numbers may not be reflected in this version. For the definitive version of this publication, please refer to the published source. You are advised to consult the publisher's version if you wish to cite this paper.

This version is being made available in accordance with publisher policies. See <http://orca.cf.ac.uk/policies.html> for usage policies. Copyright and moral rights for publications made available in ORCA are retained by the copyright holders.



1 **TITLE PAGE**

2

3 **Title: Predicting the cumulative chance of live birth over multiple complete cycles of in vitro**
4 **fertilisation: an external validation study**

5

6 **Running title:** Validation of an IVF model predicting live birth

7

8 **Authors:**

9 J.A. Leijdekkers^{1,*}, M.J.C. Eijkemans², T.C. van Tilborg¹, S.C. Oudshoorn¹, D.J. McLernon³, S.

10 Bhattacharya⁴, B.W.J. Mol⁵, F.J.M. Broekmans¹, H.L. Torrance¹, on behalf of the OPTIMIST group.

11

12 ¹Department of Reproductive Medicine and Gynaecology, University Medical Centre Utrecht, Utrecht

13 University, PO box 85500, 3508 GA Utrecht, The Netherlands. ²Julius Centre for Health Sciences and

14 Primary Care, University Medical Centre Utrecht, Utrecht University, PO box 85500, 3508GA Utrecht,

15 The Netherlands. ³Institute of Applied Health Sciences, Medical Statistics Team, University of Aberdeen,

16 Aberdeen AB25 2ZD, UK. ⁴School of Medicine, College of Biomedical and Life Sciences, Cardiff

17 University School of Medicine, Cardiff CF14 4XN, UK ⁵Department of Obstetrics and Gynaecology,

18 Monash University, VIC 3800 Clayton, Australia.

19

20 ***Correspondence address:** J.A. Leijdekkers, Department of Reproductive Medicine and Gynaecology,

21 University Medical Centre Utrecht, Utrecht University, PO box 85500, 3508 GA Utrecht, The

22 Netherlands. E-mail: j.a.leijdekkers@umcutrecht.nl

23

24 **ABSTRACT**

25

26 *Study question*

27 Are the published pre-treatment and post-treatment McLernon models, predicting cumulative live birth
28 rates (LBR) over multiple complete IVF cycles, valid in a different context?

29

30 *Summary answer*

31 With minor recalibration of the pre-treatment model, both McLernon models accurately predict
32 cumulative LBR in a different geographical context and a more recent time period.

33

34 *What is known already*

35 Previous IVF prediction models have estimated the chance of a live birth after a single fresh embryo
36 transfer, thereby excluding the important contribution of embryo cryopreservation and subsequent IVF
37 cycles to cumulative LBR. In contrast, the recently developed McLernon models predict the cumulative
38 chance of a live birth over multiple complete IVF cycles at two certain time points: a) before initiating
39 treatment using baseline characteristics (pre-treatment model) and b) after the first IVF cycle adding
40 treatment related information to update predictions (post-treatment model). Before implementation of
41 these models in clinical practice, their predictive performance needs to be validated in an independent
42 cohort.

43

44 *Study design, size, duration*

45 External validation study in an independent prospective cohort of 1515 Dutch women who participated in
46 the OPTIMIST study (NTR2657) and underwent their first IVF treatment between 2011 and 2014.
47 Participants underwent a total of 2881 complete treatment cycles, with a complete cycle defined as all
48 fresh and frozen thawed embryo transfers resulting from one episode of ovarian stimulation. The follow

49 up duration was 18 months after inclusion, and the primary outcome was ongoing pregnancy leading to
50 live birth.

51

52 *Participants/materials, setting, methods*

53 Model performance was externally validated up to three complete treatment cycles, using the linear
54 predictor as described by McLernon et al. to calculate the probability of live birth. Discrimination was
55 expressed by the c-statistic and calibration was depicted graphically in a calibration plot. In contrast to the
56 original model development cohort, anti-Müllerian hormone (AMH), antral follicle count (AFC) and body
57 weight were available in the OPTIMIST cohort, and evaluated as potential additional predictors for model
58 improvement.

59

60 *Main results and the role of chance*

61 Applying the McLernon models to the OPTIMIST cohort, the c-statistic of the *pre-treatment model* was
62 0.62 (95% confidence interval (CI) 0.59-0.64) and of the *post-treatment model* 0.71 (95% CI 0.69-0.74).
63 The calibration plot of the *pre-treatment model* indicated slight overestimation of the cumulative LBR. To
64 improve calibration, the *pre-treatment model* was recalibrated by subtracting 0.35 from the intercept. The
65 *post-treatment model* calibration plot revealed accurate cumulative LBR predictions. After addition of
66 AMH, AFC and body weight to the McLernon models, the c-statistic of the *updated pre-treatment model*
67 improved slightly to 0.66 (95% CI 0.64-0.68), and of the *updated post-treatment model* remained at the
68 previous level of 0.71 (95% CI 0.69-0.73).

69 Using the *recalibrated pre-treatment model*, a woman aged 30 years with two years of primary infertility
70 who starts ICSI treatment for male factor infertility has a chance of 40% of a live birth from the first
71 complete cycle, increasing to 72% over three complete cycles. If this woman weighs 70 kilograms, has an
72 AMH of 1.5 ng/mL and an AFC of 10 measured at the beginning of her treatment, the *updated pre-*
73 *treatment model* revises the estimated chance of a live birth to 30% in the first complete cycle and 59%
74 over three complete cycles. If this woman then has 5 retrieved oocytes, no embryos cryopreserved and a

75 single fresh cleavage stage embryo transfer in her first ICSI cycle, the *post-treatment model* estimates the
76 chances of a live birth at 28% and 58%, respectively.

77

78 *Limitations, reasons for caution*

79 Two randomised controlled trials (RCT) evaluating the effectiveness of gonadotropin dose
80 individualisation on basis of the AFC were nested within the OPTIMIST study. The strict dosing
81 regimens, the RCT in- and exclusion criteria and the limited follow up time of 18 months might have
82 influenced model performance in this independent cohort. Also, consistent with the original development
83 study, external validation was performed using the optimistic assumption that the cumulative LBR in
84 couples that discontinue treatment without a live birth would have been equal to that of couples who
85 continue treatment.

86

87 *Wider implications of the findings*

88 After national recalibration to account for geographical differences in IVF/ICSI treatment, the McLernon
89 prediction models can be introduced as new counselling tools in clinical practice to inform patients and to
90 complement clinical reasoning. These models are the first to offer an objective and personalised estimate
91 of the cumulative probability of live birth over multiple complete IVF cycles.

92

93 *Study funding/competing interest(s):*

94 No external funds were obtained for this study. M.J.C.E., D.J.M. and S.B. have nothing to disclose. J.A.L.,
95 S.C.O, T.C.v.T. and H.LT. received an unrestricted personal grant from Merck BV. B.W.M. is supported
96 by a NHMRC Practitioner Fellowship (GNT1082548) and reports consultancy for ObsEva, Merck and
97 Guerbet. F.J.M.B. receives monetary compensation as a member of the external advisory board for Merck
98 BV (the Netherlands) and Ferring pharmaceuticals BV (the Netherlands), for consultancy work for Gedeon
99 Richter (Belgium) and Roche Diagnostics on automated AMH assay development, and for a research
100 cooperation with Ansh Labs (USA).

101

102 *Trial registration number*

103 Not applicable

104

105 **KEYWORDS**

106 Prediction model, external validation, live birth, IVF/ICSI, infertility, cumulative live birth, personalised,

107 counselling, prognostic research

108 **Introduction**

109 Infertility is defined as the failure to conceive within 12 months of regular unprotected intercourse, and
110 affects approximately one in six couples (Oakley *et al.*, 2008; Zegers-Hochschild *et al.*, 2017). The
111 majority of infertile couples seek fertility care, and many of those with prolonged unresolved infertility
112 will be treated with ART regardless of cause (Boivin *et al.*, 2007; Datta *et al.*, 2016). IVF and ICSI are
113 both widely used techniques for couples with infertility. Globally more than 1.6 million annual cycles of
114 IVF/ICSI are performed and while success rates have increased over time (Dyer *et al.*, 2016; McLernon *et*
115 *al.*, 2016), this treatment is still not effective for all infertile couples, with live birth rates (LBR) at around
116 25-30% per treatment cycle (Malizia *et al.*, 2009; McLernon *et al.*, 2016; De Neubourg *et al.*, 2016).
117 Since IVF/ICSI is expensive and carries several risks, the probability of a live born child should be
118 weighed against the risks and costs of this treatment.

119 Several prognostic models have been developed to objectively estimate the probability of a live birth after
120 IVF/ICSI treatment (Leushuis *et al.*, 2009; van Loendersloot *et al.*, 2014). It is known that prediction
121 models often perform optimistically in their development sample, even after correction by internal
122 validation. This is caused by overfitting, which occurs when the model corresponds too closely to the
123 development data due to the inclusion of too many predictors (Moons, Kengne, Woodward, *et al.*, 2012).
124 External validation in an independent cohort of women is thus essential to examine the performance and
125 generalisability of the prediction model (Altman *et al.*, 2009; Harrell *et al.*, 1996). Unfortunately, most of
126 the currently available models that predict the chance of a live birth after IVF/ICSI treatment have never
127 been externally validated (Leushuis *et al.*, 2009; van Loendersloot *et al.*, 2014). Also, the majority of
128 these models predict the probability of a live birth after a single fresh embryo transfer, excluding the
129 important contribution of embryo cryopreservation and subsequent treatment cycles to LBR. This limits
130 their potential as counselling tools for couples and clinicians, especially considering the increased use and
131 improved techniques of embryo cryopreservation and frozen thawed embryo transfer cycles in recent
132 years (Wong *et al.*, 2014).

133 Three of the largest model development studies for prediction of live birth after IVF and/or ICSI
134 treatment used data from the Human Fertilisation and Embryology Authority (HFEA) database in the UK
135 (McLernon *et al.*, 2016; Nelson and Lawlor, 2011; Templeton *et al.*, 1996). Treatment and outcome data
136 from all licenced fertility clinics within the UK have been recorded in this database since 1992. The two
137 models developed by Templeton *et al.* and Nelson *et al.* were both externally validated, and their
138 predictive performance was compared to one another in several studies (Arvis *et al.*, 2012; van
139 Loendersloot *et al.*, 2011; Smeenk *et al.*, 2000; Smith *et al.*, 2015; te Velde *et al.*, 2014). Although these
140 models have been recommended in previous studies and used internationally to predict live birth after
141 IVF and ICSI (Leushuis *et al.*, 2009; Smith *et al.*, 2015; te Velde *et al.*, 2014), neither model predicts
142 cumulative LBR over multiple IVF/ICSI treatment cycles including frozen thawed embryo transfer
143 cycles.

144 Recently, a new model was developed by McLernon *et al.* using the HFEA database (McLernon *et al.*,
145 2016). This model is the first to provide an individualised estimate of the cumulative chance of a live
146 birth over multiple complete cycles of IVF/ICSI, with a complete cycle defined as all fresh and frozen
147 thawed embryo transfers resulting from one episode of ovarian stimulation. For model development, data
148 from 113 873 women and 184 269 complete cycles between 1999 and 2009 were used. Internal validation
149 of the model showed promising results, however evaluation of the predictive performance of the model in
150 a different geographical context using more contemporary data has yet to be performed. Additionally, a
151 number of potential key predictors, such as measures for ovarian reserve and female body weight, were
152 unavailable in the HFEA database and could not be included in the original model (McLernon *et al.*,
153 2016).

154 The main objective of the current study was therefore to perform geographical and temporal validation of
155 the new HFEA model by using recent data from a different country. We also wanted to determine whether
156 inclusion of additional parameters, such as female body weight and ovarian reserve test results i.e. antral

157 follicle count (AFC) and anti-Müllerian hormone (AMH), could improve the predictive performance of
158 the model.

159 **Materials and methods**

160 *Data sources*

161 External validation was performed on data from the OPTIMIST study (van Tilborg, Oudshoorn, *et al.*,
162 2017). This multicentre prospective cohort study included 1515 women from 25 infertility centres in the
163 Netherlands between May 2011 and May 2014. Participants were younger than 44 years of age, had
164 regular menstrual cycles and no significant uterine or ovarian abnormalities on transvaginal ultrasound.
165 Women with polycystic ovarian syndrome, metabolic or endocrine abnormalities or undergoing oocyte
166 donation were excluded. All participants were included before their first IVF/ICSI cycle, or the first cycle
167 after a previous live birth. The primary outcome was ongoing pregnancy, achieved within 18 months of
168 follow up, and resulting in live birth. Ethical approval for the OPTIMIST study was obtained from the
169 Institutional Review Board of the University Medical Centre Utrecht (MEC 10-273), and all participants
170 provided written informed consent. A more detailed description of study procedures and results were
171 reported previously (Oudshoorn *et al.*, 2017; van Tilborg *et al.*, 2012; van Tilborg, Oudshoorn, *et al.*,
172 2017; van Tilborg, Torrance, *et al.*, 2017).

173 *McLernon model*

174 The McLernon model consists of two clinical prediction models to estimate the individualised cumulative
175 chance of a live birth over a maximum of six complete treatment cycles. **Before initiating treatment**, the
176 *pre-treatment model* predicts the probability of a live birth from both fresh and frozen thawed embryo
177 transfers based on couple characteristics and the use of IVF or ICSI. Included predictors are: female age
178 (years), duration of infertility (years), previous pregnancy, causes of infertility (tubal factor, anovulation,
179 male factor, unexplained infertility), type of treatment (IVF or ICSI) and treatment year (see
180 Supplementary Text 1).

181 **After the first fresh treatment cycle**, treatment specific characteristics from this cycle are added in the
182 *post-treatment model* to update the predicted probability. Added predictors are: number of oocytes,

183 cryopreservation of embryos, and the number and stage of embryos at the first fresh embryo transfer
184 (single, double or triple embryo transfer; blastocyst or cleavage stage). All causes of infertility are
185 excluded as predictors in the post-treatment model, except for tubal factor (see Supplementary Text 2).
186 For women with zero oocytes collected in the first cycle, a separate post-treatment model is available.

187 To predict the probability of a live birth in the i th cycle, assuming no live birth occurred in the previous
188 cycle(s), complete cycle number is included in both models as a discrete time variable. A complete cycle
189 includes all fresh and frozen thawed embryo transfers resulting from one episode of ovarian stimulation.
190 With the predicted probability of a live birth per subsequent complete cycle, the cumulative probability of
191 a live birth can be calculated up to six complete cycles (see Supplementary Text 1 and 2).

192 *Statistical analysis*

193 Nine predictor variables had missing values (Table I). The proportion of missing values was low (<
194 2.5%), except for AMH (11.2%). During the OPTIMIST study, blood sampling was performed on the day
195 of randomisation. Logistic issues prevented blood sampling in some cases, thus compromising the ability
196 to undertake post-hoc measurements of AMH in the total population. As the reasons for missing values
197 were considered to be unrelated to the AMH value itself or the measurement, these were defined as
198 missing (completely) at random.

199 Multiple imputation was applied for predictors with missing values in the OPTIMIST database (Sterne *et*
200 *al.*, 2009). In this process 10 imputed datasets were created using a multivariate imputation by chained
201 equations (MICE) algorithm (van Buuren and Groothuis-Oudshoorn, 2011). Predicted probabilities for a
202 live birth were calculated on each imputed dataset, using the predictors and parameter-estimates of both
203 the pre-treatment model as well as the post-treatment model as described by McLernon *et al* 2016
204 (McLernon *et al.*, 2016). In accordance with the original models, the variables female age, treatment year
205 and number of oocytes were treated with restricted cubic splines in the validation process. The separate
206 post-treatment model for women with zero oocytes collected in the first treatment cycle was not validated

207 in this study, as the number of women for this analysis was too low in the OPTIMIST database.
208 Cumulative probabilities were calculated up to three complete IVF/ICSI cycles, as most couples in the
209 Netherlands only have three treatment cycles due to the current reimbursement policy. Also, the
210 OPTIMIST follow up period was 18 months, reducing the number of women with more than three
211 treatment cycles. The validation process was performed ten times on each of the imputed datasets and
212 separate results were pooled using Rubin's rules (Rubin, 2004).

213 The predictive performance of the McLernon models was evaluated in terms of discrimination and
214 calibration. Discrimination quantifies the ability of a model to correctly differentiate between subjects
215 with an event and subjects without an event (Moons, Kengne, Woodward, *et al.*, 2012). In the context of
216 fertility treatment, it is the ability of the models to distinguish between women with a live birth and
217 women without a live birth after IVF/ICSI treatment. It is expressed by the c-statistic or the area under the
218 receiver operating curve (AUROC), which ranges between 0.5 and 1. A c-statistic of 1 indicates perfect
219 discrimination, whereas a c-statistic of 0.5 represents a model with no discrimination at all. In this study,
220 the c-statistic (and 95% CI) was calculated using the method suggested by Harrell *et al.* (Harrell *et al.*,
221 1996).

222 Calibration describes the degree of agreement between predicted probabilities and observed outcomes
223 (Moons, Kengne, Woodward, *et al.*, 2012), in this context the predicted probability of a live birth and the
224 observed LBR. Calibration can be assessed graphically by forming subgroups of patients determined by
225 ranges of predicted probabilities, and then plotting the observed proportion of events against the mean
226 predicted probability within these subgroups. When perfect calibration is present, the plot shows a
227 diagonal line with a slope of one and an intercept of zero. In the current study, five equal subgroups of
228 patients were formed. This was based on the sample size of the OPTIMIST cohort and the related
229 precision of the point estimates in the calibration plot. Within these subgroups, the Kaplan Meier
230 estimates of the observed cumulative LBR over three complete treatment cycles were plotted against the
231 mean predicted probability of cumulative live birth. A smoothed line was then added in this plot using the

232 proportional hazard regression approach described by Harrell et al (Harrell *et al.*, 1996). In addition to
233 this, a systematic difference in the predicted and observed LBR was assessed by using calibration-in-the-
234 large (Steyerberg, 2009), and the intercept of the prediction models was adjusted in case a systematic
235 over- or underestimation was present.

236 *Updating the models*

237 Following the external validation of the models, the additional value of updating the McLernon models
238 with pre-specified new biomarkers was evaluated. AMH (ng/mL), AFC (2-10 mm) and body weight (kg)
239 were added to the pre-treatment and post-treatment model in a multivariable logistic regression analysis,
240 in which the linear predictor of the McLernon model was entered as a fixed variable. The final model was
241 established using a manual backward selection process. Predictors were eliminated from the model
242 according to the Akaike Information Criterion (AIC) (Akaike, 1974).

243 The predictive performance of the new updated models was evaluated by calculating the c-statistic (and
244 95% CI). To assess for overfitting, internal validation was performed by bootstrapping (Steyerberg,
245 2009). Two hundred bootstrap samples, all of which were of the same size as the original validation
246 sample, were created by random sampling with replacement (Harrell, 2001; Steyerberg, 2009). In each
247 bootstrap sample, a new model was fitted with the same predictors as the updated models. The c-statistic
248 was calculated for each of the 200 sample derived models, in both the bootstrap sample as well as the
249 original validation cohort. The difference between these two c-statistics was calculated for each of the 200
250 sample derived models, and averaged to give the optimism estimate. This was subtracted from the
251 original c-statistic to obtain the optimism corrected c-statistic for the updated models.

252 All statistical analyses were performed using R for Windows (version 3.3.2; R Foundation for Statistical
253 Computing, Vienna, Austria).

254 **Results**

255 Of the 1515 women included in the OPTIMIST study, four were excluded in the current study as they
256 never started IVF/ICSI treatment. A total of 2881 IVF/ICSI cycles were performed over a period of 18
257 months of follow up. Table I shows the patient and first cycle treatment characteristics of the OPTIMIST
258 cohort (validation sample) and the HFEA cohort (development sample). Women included in the
259 validation sample were about the same age as women in the development sample, but had a shorter
260 average duration of infertility. The causes of infertility showed a similar distribution across both samples,
261 with the exception of anovulation which rendered women ineligible for the OPTIMIST study. The
262 treatment characteristics showed that embryo cryopreservation was more frequently performed after the
263 first IVF/ICSI cycle in the validation sample and that these women most often had a cleavage stage single
264 embryo transfer in the first fresh cycle, whereas women in the development sample most often had a
265 cleavage stage double embryo transfer. No formal assessment was performed for the differences and
266 similarities between the cohorts, as a description rather than a p-value is considered to be useful for
267 interpretation of the models' performance in this external validation study.

268 The flowchart in Figure 1 shows the number of women in the OPTIMIST and HFEA cohorts who started
269 a treatment cycle, had a live birth or discontinued treatment without having a live birth. The LBR per
270 cycle was similar in both cohorts for the first, second and fourth treatment cycle. In the third cycle the
271 LBR was slightly higher in the OPTIMIST cohort compared to the HFEA cohort. As few women in the
272 OPTIMIST cohort received a fifth or sixth cycle, LBR in these cycles could not be compared. The
273 proportion of women without a live birth that continued treatment was higher after the first and second
274 cycle in the OPTIMIST cohort as compared to the HFEA cohort. After the third cycle, the proportion
275 continuing treatment in the OPTIMIST cohort decreased, while it remained constant in the HFEA cohort.
276 At the end of follow up, 52% of the women in the OPTIMIST study had a treatment related live birth. The
277 overall LBR of the HFEA cohort was 43% over six complete IVF/ICSI cycles.

278 As mentioned previously, external validation of the McLernon models was performed up to three
279 complete treatment cycles, and therefore the fourth, fifth and sixth complete treatment cycle in the
280 OPTIMIST dataset (n=102 complete treatment cycles, n= 15 live births) were excluded from further
281 analysis. Also, for the post-treatment model validation, women with zero oocytes collected in the first
282 treatment cycle were excluded (n= 226 women, n = 526 complete treatment cycles, n= 82 live births) as a
283 separate model was developed for this group of women by McLernon et al (McLernon *et al.*, 2016). Due
284 to the small numbers, this separate model could not be validated in this study.

285 *Discrimination and calibration*

286 In the validation sample, the pooled c-statistic for the pre-treatment model was 0.62 (95% CI 0.59-0.64)
287 and for the post-treatment model 0.71 (95% CI 0.69-0.74). Figure 2a and 3 show the calibration plots for
288 both original models, depicting the correlation between the observed and predicted cumulative LBR. The
289 pre-treatment calibration plot had an intercept of -0.23 (95% CI -0.36- -0.10) and a slope of 0.98 (95% CI
290 0.69-1.27), and the post-treatment calibration plot had an intercept of -0.01 (95% CI -0.12-0.11) and a
291 slope of 0.97 (95% CI 0.77-1.19).

292 The pre-treatment model systematically overestimated the cumulative LBR over three complete cycles for
293 women in the validation sample. This is shown by a calibration curve with most of the confidence
294 intervals under the reference line (Figure 2a), indicating significantly higher predicted probabilities than
295 observed LBR. The calibration-in-the-large analysis confirmed this systematic overestimation with an
296 intercept of -0.35. To improve calibration, the pre-treatment model was thus adjusted by subtracting 0.35
297 from the intercept of the original linear predictor, which decreased the predicted odds of a live birth by a
298 factor of 1.42 (see Supplementary Text 3). The calibration plot of the recalibrated pre-treatment model
299 showed improved accuracy of the predictions, with all confidence intervals overlapping the reference line
300 (Figure 2b). In contrast to the pre-treatment model, the post-treatment model correctly estimated the

301 cumulative LBR in the validation sample, as is shown by a calibration plot with confidence intervals
302 overlapping the reference line indicating no significant over- or underestimation (Figure 3).

303 *Updating of the models*

304 Addition of the biomarkers AMH, AFC and body weight to the pre-treatment and post-treatment model in
305 a multivariable regression analysis resulted in two new updated models. The updated pre-treatment model
306 included all three biomarkers as additional predictors for live birth. Since the relationship between both
307 AMH and AFC with the probability of live birth was non-linear, these predictors were included using
308 restricted cubic splines (see Supplementary Figure 1). The updated post-treatment model included only
309 AFC and AMH as additional predictors for live birth, of which AFC was modelled by using restricted
310 cubic splines (see Supplementary Figure 2). After internal validation of the updated models by
311 bootstrapping, the updated pre-treatment model had a corrected c-statistic of 0.66 (95% CI 0.64-0.68) and
312 the updated post-treatment model had a corrected c-statistic of 0.71 (95% CI 0.69-0.73). The addition of
313 AFC, AMH and body weight thus resulted in a slight improvement of the discriminatory capacity of the
314 pre-treatment model, while addition of AFC and AMH had no beneficial effect on the discriminative
315 performance of the post-treatment model.

316 *Examples of model predictions*

317 Figures 4, 5 and 6 show examples of model predictions as illustration for clinical application. Figure 4
318 presents predictions of the *recalibrated pre-treatment model* for couples with primary infertility caused
319 by a male factor. Cumulative probabilities of live birth are calculated up to three complete ICSI cycles,
320 and are differentiated by female age (30 or 40 years) and duration of infertility (2 years or 5 years). As is
321 shown in figure 4, age is the most important predictor in the pre-treatment model. A 30-year-old woman
322 with 2 years of infertility has a predicted probability of a live birth of 0.40 in the first ICSI cycle,
323 increasing to 0.72 over three complete cycles. For a 40-year-old woman with 2 years of infertility, these
324 probabilities are 0.15 and 0.32 respectively.

325 Figure 5 shows predictions of the *updated pre-treatment model*, with AMH, AFC and body weight as new
326 predictors in the model. Predictions are presented for couples with two years of primary infertility caused
327 by a male factor, and differentiation is based on female age (30 or 40 years), AMH (2.0 or 0.5 ng/mL) and
328 AFC (15 or 7). In all scenarios the female body weight is 70 kilograms. A 30-year-old woman with an
329 average ovarian reserve at the start of her first treatment – indicated by an AMH of 2.0 ng/mL and an
330 AFC of 15 – has a predicted probability of a live birth of 0.37 in the first cycle and 0.69 over three cycles
331 (0.17 and 0.37 for a 40-year-old woman). If this woman has a reduced ovarian reserve – indicated by an
332 AMH of 0.5 ng/mL and an AFC of 7 – the predicted probabilities decrease to 0.19 and 0.42, respectively
333 (0.08 and 0.18 for a 40-year-old woman).

334 Figure 6 shows predictions of the *post-treatment model*, which revises the predicted probabilities of the
335 pre-treatment models by adding information of the first treatment cycle. Predictions are calculated for
336 women with two years of primary, non-tubal infertility and are differentiated by female age (30 or 40
337 years), number of oocytes (10 or 5) and embryo cryopreservation (yes or no). In all scenarios the woman
338 received a cleavage stage single embryo transfer. The predicted probabilities of a live birth for women
339 with a favourable prognosis – aged 30-years, 10 oocytes retrieved and cryopreserved embryos – is 0.49 in
340 the first ICSI cycle, increasing to 0.83 over 3 complete cycles. In contrast, for women with a poorer
341 prognosis – aged 40 years, 5 oocytes retrieved and no embryos cryopreserved – the predicted probabilities
342 are 0.11 and 0.26, respectively.

343 **Discussion**

344 *Main findings*

345 This external validation study of the McLernon pre-treatment and post-treatment model found that, after
346 minor recalibration of the intercept of the pre-treatment model, both models accurately predict the
347 cumulative probability of live birth up to three complete IVF/ICSI cycles in a more contemporary cohort
348 in another country. The discriminatory capacity of the pre-treatment model in an external cohort was
349 limited, whereas the post-treatment model had a fair ability to discriminate between couples with and
350 without a live birth after treatment.

351 *Strengths*

352 This study focuses on the external validation of an IVF prediction model, which is an essential but
353 frequently overlooked step before implementation of a prediction model in clinical practice (Altman *et*
354 *al.*, 2009). In contrast to redeveloping new models for the same outcome, external validation and updating
355 of existing models prevents the loss of scientific information by combining the information captured in
356 the original model with information of a new patient cohort (Moons, Kengne, Grobbee, *et al.*, 2012).

357 Embryo cryopreservation has become an important part of IVF/ICSI treatment, and most couples have
358 more than just one complete treatment cycle (Wong *et al.*, 2014). Unlike previous prediction models
359 (Leushuis *et al.*, 2009; van Loendersloot *et al.*, 2014), the McLernon models provide a more useful
360 estimate of cumulative treatment success. As such, the validation of these models represents a significant
361 step forward in creating a clinically useful tool to manage expectations and to inform decision making
362 around IVF.

363 This study benefits from the prospective design of the OPTIMIST study, which has ensured reliable data
364 collection, with relatively low numbers of missing values and a low risk of selection bias. The multicentre
365 design resulted in a highly representable cohort for Dutch fertility care. And although it is known that the

366 IVF/ICSI success rates vary between fertility centres, the inclusion of multiple centres will increase the
367 generalisability and applicability of the external validation of the McLernon models within the
368 Netherlands.

369 Furthermore, the external validation was performed on data collected in a recent time period (2011-2014).
370 Due to changing patient populations, new treatment protocols, improving technologies and increasing
371 success rates over time, prediction models in reproduction medicine have no static form and should be
372 regularly updated to optimally reflect the latest circumstances in which they are used (Altman *et al.*,
373 2009). As the McLernon models were developed on data collected between 1999 and 2009, data of the
374 more recently performed OPTIMIST study were helpful to investigate if model performance was still
375 accurate in current practice.

376 *Weaknesses*

377 This study has a number of limitations. First, the external validation involved data from a prospective
378 cohort study within which two randomised controlled trials were embedded evaluating the effectiveness
379 of individualised doses of gonadotropins based on AFC. Strict dosing regimens might have affected some
380 treatment outcomes, such as cancellation rates and number of oocytes, thus influencing the predictive
381 capacity of the models in the validation sample. However, as the OPTIMIST study found no difference
382 between the dosing regimens on cumulative live birth rates, the impact on model performance is likely to
383 be minimal.

384 Second, the OPTIMIST study used strict eligibility criteria. Therefore, the validation sample does not
385 fully represent the diversity of the patient population initiating IVF/ICSI treatment in the Netherlands. As
386 none of the women in the validation sample were anovulatory, external validation of the models was only
387 performed for an ovulatory population. This limits the generalisability of the models to some extent, as
388 the original McLernon models were developed in a population which also included anovulatory women.
389 Also, it could have had some impact on model performance. However, since anovulation had only a small

390 predictive value in the pre-treatment model, and the majority of couples underwent IVF/ICSI for other
391 indications, a large impact on model performance is unlikely.

392 Third, the OPTIMIST study had a follow up period of 18 months, leading to small numbers of women
393 with more than three complete treatment cycles. Model performance could therefore only be reliably
394 validated up to three complete cycles. However, most couples in the Netherlands complete a maximum of
395 three treatment cycles which is partly due to the national reimbursement policy, but also by the high rates
396 of embryo cryopreservation, increasing the number of embryo transfers and LBR per cycle. Therefore,
397 model validation up to three complete cycles has particular clinical relevance for current Dutch fertility
398 care.

399 Last, the original McLernon prediction models were developed on linked cycle data, which were then
400 used to estimate cumulative pregnancy chances. Therefore, these models used the optimistic assumption
401 that the cumulative LBR in couples who discontinue IVF treatment without a live birth would have been
402 equal to that of couples who continue further treatment cycles, after correction of predictor effects. This
403 assumption tends to lead to overestimation of the cumulative LBR, as women with a low prognosis of
404 achieving a live birth are generally more likely to discontinue treatment (Brandes *et al.*, 2009; Olivius *et*
405 *al.*, 2004). Since the reasons for treatment withdrawal were unknown in the current external validation
406 study, a similar method was used that probably resulted in some degree of overestimation of the
407 cumulative LBR in the validation cohort. However, as the original McLernon models were developed
408 with this approach, and the predictions for cumulative LBR over multiple complete cycles were
409 considered to be clinically more relevant than per cycle predictions, we feel that the current method is the
410 best option for the external validation of the McLernon models.

411 *Explanation of findings*

412 The discriminatory capacity of the pre-treatment model was markedly lower in the validation sample than
413 in the development sample. In the development study, a c-statistic of 0.73 (95% CI 0.72-0.74) was

414 reported, whereas the present study found a c-statistic of 0.62 (95% CI 0.59-0.64). For the post-treatment
415 model, the discriminatory performance in the validation sample was comparable to that in the
416 development sample, with a c-statistic of 0.71 (95% CI 0.69-0.74) and 0.72 (95% CI 0.71-0.73)
417 respectively (McLernon *et al.*, 2016). As it is known that prediction models tend to perform too
418 optimistically in the development dataset due to overfitting, some reduction in model performance is to be
419 expected during external validation due to the differences between samples (Altman *et al.*, 2009; Moons,
420 Kengne, Woodward, *et al.*, 2012). This, to some extent, also explains the lower overall performance of
421 the pre-treatment model. The comparable performance of the post-treatment model in both samples
422 indicates that the treatment related variables that were added to this model (number of oocytes,
423 cryopreservation of embryos, and the number and stage of embryos) are important predictors for live birth
424 after treatment.

425 Other than the influence of overfitting, some key differences between the Dutch and UK healthcare
426 systems may also have affected the models' performance in this external validation study. An important
427 factor is the reimbursement policy for fertility treatment. All Dutch infertile couples are insured for a
428 minimum of three complete IVF/ICSI cycles. In contrast, most couples in the UK receive no standard
429 funding for ART (Berg Brigham *et al.*, 2013). Since IVF/ICSI treatment is expensive, this induces
430 discrepancies in the patient population initiating and continuing treatment between the two study samples
431 (Rajkhowa *et al.*, 2006). As can be seen in the baseline table (Table I) and flowchart (Figure 1), couples
432 in the UK had a longer average duration of infertility before starting treatment and were more likely to
433 discontinue treatment after the first and second cycles than couples in the Netherlands. Also, the decrease
434 in LBR is more evident in the UK than in the Netherlands over the first three cycles, which suggests that
435 differences exist in both reasons for discontinuation as well as prognostic profiles of women
436 discontinuing treatment in the two countries. These phenomena are, in part, financially driven, and could
437 partially explain the difference in predictive ability of the UK models in the Dutch cohort.

438 Furthermore, despite the fact that the infertility guidelines of both countries include similar approaches
439 for treatment of infertile couples, there are important variations in treatment characteristics between the
440 two study samples (Dutch Society of Obstetrics and Gynaecology (NVOG), 2010; National Institute for
441 Health and Care Excellence (NICE), 2013). Some of these differences are mainly due to changes in
442 clinical practice over time. As is shown by the baseline table (Table 1), women in the more recent Dutch
443 cohort (2011-2014) generally had a single embryo transfer in their first fresh treatment cycle, whereas
444 women in the earlier UK cohort (1999-2009) most often had a double embryo transfer. Also, embryo
445 cryopreservation was performed in over half of the Dutch women as compared to only a quarter of the
446 women in the UK. Other differences are explained by variation in treatment protocols between
447 geographic locations. For one, no blastocyst stage embryos transfers were performed in the Netherlands in
448 contrast to the proportion of blastocyst stage embryo transfers in the UK of more than 10%. Also, Dutch
449 women more frequently had no embryo available for transfer after their first treatment cycle, which is
450 most likely caused by strict cancellation criteria particularly for hyper response. These differences in
451 treatment characteristics suggest that the development sample does not fully reflect clinical practice in a
452 more recent time period and in a different geographic context. As cumulative LBR are substantially
453 affected by the variation in treatment characteristics (Glujovsky *et al.*, 2016; Pandian *et al.*, 2013; Wong
454 *et al.*, 2014), this could explain part of the different performance of the pre-treatment model in the
455 validation sample . The stable performance of the post-treatment model, which includes embryo stage and
456 embryo cryopreservation as important predictors, seems to confirm the impact of the variation in these
457 variables on model performance.

458 The addition of measures of ovarian reserve, i.e. AMH and AFC, and body weight to the McLernon
459 prediction models revealed only a marginal improvement of model performance in the OPTIMIST
460 dataset. The additional value of these tests can therefore be questioned, especially in view of the extra
461 costs and physical burden on the patient. Female age is one of the most important predictors in the
462 McLernon models (McLernon *et al.*, 2016). As female age is correlated with the ovarian reserve, adding

463 AMH and AFC provides limited new information to the prediction models. This is in line with previous
464 studies that showed that ovarian reserve tests have no added value to the use of female age alone in the
465 prediction of ongoing pregnancy after treatment (Broer *et al.*, 2013). Other potential predictors for live
466 birth, such as ethnicity, smoking status and alcohol intake, were not included in this update of the
467 McLernon model (Dhillon *et al.*, 2015; Rossi *et al.*, 2011; Waylen *et al.*, 2009). The additional value of
468 these variables for model performance was considered uncertain, as the reporting is remarkably subjective
469 and/or often incomplete (Liber and Warner, 2018; Stockwell *et al.*, 2016).

470 *Clinical implications*

471 Discrimination and calibration have been recognized as measures to evaluate the performance of
472 prediction models (Altman *et al.*, 2009; Steyerberg, 2009). However, the discriminative ability at the
473 binary level of most prediction models in reproductive medicine, as expressed by the c-statistic, is
474 considerably low (Leushuis *et al.*, 2009). As at the moment of prediction the outcome of pregnancy has
475 not yet occurred, the c-statistic is determined using the calculated probability of pregnancy. The
476 maximum value of the c-statistic depends on the variability of these calculated probabilities in the
477 infertile population. Since infertility is a complex and multifactorial health problem and due to the
478 absence of strong predictors for live birth – particularly pre-treatment – , the probability distribution in
479 infertile couples that have a live birth has a considerable overlap with the distribution of those without a
480 live birth. Therefore the maximum c-statistic can be expected to be low (Cook, 2007; Coppus *et al.*,
481 2009), as is seen in the external validation of the pre-treatment model. However, this does not necessarily
482 imply that such prediction models have limited use in clinical practice. Models with reliable predictions
483 and a clinically useful distribution of probabilities for achieving a live birth, as assessed by calibration,
484 can still support patients and clinicians in clinical decision making around infertility treatment (Coppus *et*
485 *al.*, 2009).

486 As the calibration plots of both the recalibrated pre-treatment model and the post-treatment model
487 indicate accurate predictions with a useful range of prognoses, these models can be used within the
488 Netherlands as counselling tools to complement clinical reasoning at two certain time points. Before
489 initiating treatment, the recalibrated pre-treatment model offers couples and clinicians a personalised and
490 objective estimate of success over multiple complete treatment cycles. And after the first fresh embryo
491 transfer, the post-treatment model provides a revised estimate using treatment related information to
492 personalize the predictions even more. Despite the applicability of the models as counselling tools to
493 inform patients about their prognosis, the McLernon models should not yet be used for decisions on
494 whether or not to withhold fertility treatment. The impact of such model-based decisions on cost-benefit
495 outcomes should be investigated first and proven to be beneficial. To implement the McLernon models as
496 counselling tools in other countries as well, national recalibration is recommended to account for
497 geographical differences in IVF/ICSI treatment.

498 The McLernon models were converted into an online calculator to facilitate the use of the models in
499 clinical practice (<https://w3.abdn.ac.uk/clsm/opis>). As the original pre-treatment model overestimates
500 cumulative LBR for couples in the Netherlands, conversion of the recalibrated pre-treatment model into a
501 new online calculator is needed for implementation in Dutch clinical practice. This tailored online
502 calculator can then provide accurate and up to date predictions for couples and clinicians in the
503 Netherlands. Ultimately, the online calculator will be offered for implementation on the websites of the
504 Dutch Patient Association for people with fertility problems 'Freya' and the Dutch Association of
505 Obstetrics and Gynaecology (NVOG) to increase the accessibility of the models.

506 *Research implications*

507 Following this external validation study, future studies could focus on the impact of introducing the
508 McLernon prediction models in clinical practice, and assess changes in patient and clinicians' behaviour
509 and its effects on LBR and cost-effectiveness.

510 In conclusion, after minor recalibration of the pre-treatment model, the McLernon models have proven to
511 be valid in predicting the chance of cumulative live birth after multiple complete treatment cycles in
512 another geographical context and in a more recent time period. Updating the models with AMH, AFC and
513 body weight revealed only a marginal improvement of predictive performance. Following national
514 recalibration, implementation of the McLernon models as counselling tools in clinical practice will
515 provide infertile couples and clinicians with objective and personalized estimates of success over multiple
516 complete IVF/ICSI cycles.

517

518 **Acknowledgements**

519 We would like to thank the women who participated in the OPTIMIST study and the staff of the
520 participating hospitals for their contributions to the OPTIMIST study.

521 **Author's roles**

522 T.C.v.T. and S.C.O and all other members from the OPTIMIST study group collected the data. D.J.M.,
523 S.B., F.J.M.B. and H.L.T were involved in study conception and study design. J.A.L. and M. J. C. E.
524 performed the statistical analysis. J.A.L. drafted the manuscript. J.A.L., M.J.C.E., F.J.M.B. B.W.M.,
525 H.L.T interpreted the data. All authors participated to the discussion of the findings and revised the
526 manuscript.

527 **Funding**

528 No external funding was obtained for this study.

529 **Conflict of interest**

530 M.J.C.E., D.J.M. and S.B. have nothing to disclose. J.A.L, S.C.O, T.C.v.T. and H.L.T. received an
531 unrestricted personal grant from Merck BV. B.W.M. is supported by a NHMRC Practitioner Fellowship
532 (GNT1082548) and reports consultancy for ObsEva, Merck and Guerbet. F.J.M.B. receives monetary
533 compensation as a member of the external advisory board for Merck Serono (the Netherlands) and
534 Ferring pharmaceuticals BV (the Netherlands), for consultancy work for Gedeon Richter (Belgium) and
535 Roche Diagnostics on automated AMH assay development, and for a research cooperation with Ansh
536 Labs (USA).

537 **References**

- 538 Akaike H. A new look at the statistical model identification. *IEEE Trans Automat Contr* 1974;**19**:716–
539 723.
- 540 Altman DG, Vergouwe Y, Royston P, Moons KGM. Prognosis and prognostic research: validating a
541 prognostic model. *BMJ* 2009;**338**:b605.
- 542 Arvis P, Lehert P, Guivarc’h-Leveque A. Simple adaptations to the Templeton model for IVF outcome
543 prediction make it current and clinically useful. *Hum Reprod* 2012;**27**:2971–2978.
- 544 Berg Brigham K, Cadier B, Chevreur K. The diversity of regulation and public financing of IVF in
545 Europe and its impact on utilization. *Hum Reprod* 2013;**28**:666–75.
- 546 Boivin J, Bunting L, Collins JA, Nygren KG. International estimates of infertility prevalence and
547 treatment-seeking: potential need and demand for infertility medical care. *Hum Reprod*
548 2007;**22**:1506–12.
- 549 Brandes M, van der Steen JOM, Bokdam SB, Hamilton CJCM, de Bruin JP, Nelen WLDM, Kremer
550 JAM. When and why do subfertile couples discontinue their fertility care? A longitudinal cohort
551 study in a secondary care subfertility population. *Hum Reprod* 2009;**24**:3127–35.
- 552 Broer SL, van Disseldorp J, Broeze KA, Dolleman M, Opmeer BC, Bossuyt P, Eijkemans MJC, Mol B-
553 WJ, Broekmans FJM, Broer SL, *et al*. Added value of ovarian reserve testing on patient
554 characteristics in the prediction of ovarian response and ongoing pregnancy: an individual patient
555 data approach. *Hum Reprod Update* 2013;**19**:26–36.
- 556 van Buuren S, Groothuis-Oudshoorn K. MICE : Multivariate Imputation by Chained Equations in R. *J*
557 *Stat Softw* 2011;**45**:1–67.
- 558 Cook NR. Statistical Evaluation of Prognostic versus Diagnostic Models: Beyond the ROC Curve. *Clin*

- 559 *Chem* 2007;**54**:17–23.
- 560 Coppus SFPJ, van der Veen F, Opmeer BC, Mol BWJ, Bossuyt PMM. Evaluating prediction models in
561 reproductive medicine. *Hum Reprod* 2009;**24**:1774–1778.
- 562 Datta J, Palmer MJ, Tanton C, Gibson LJ, Jones KG, Macdowall W, Glasier A, Sonnenberg P, Field N,
563 Mercer CH, *et al.* Prevalence of infertility and help seeking among 15 000 women and men. *Hum*
564 *Reprod* 2016;**31**:2108–2118.
- 565 Dhillon RK, Smith PP, Malhas R, Harb HM, Gallos ID, Dowell K, Fishel S, Deeks JJ, Coomarasamy A.
566 Investigating the effect of ethnicity on IVF outcome. *Reprod Biomed Online* 2015;**31**:356–363.
- 567 Dutch Society of Obstetrics and Gynaecology (NVOG). Landelijke Netwerkrichtlijn Subfertiliteit. 2010.
- 568 Dyer S, Chambers GM, de Mouzon J, Nygren KG, Zegers-Hochschild F, Mansour R, Ishihara O, Banker
569 M, Adamson GD. International Committee for Monitoring Assisted Reproductive Technologies
570 world report: Assisted Reproductive Technology 2008, 2009 and 2010. *Hum Reprod* 2016;**31**:1588–
571 1609.
- 572 Glujovsky D, Farquhar C, Quinteiro Retamar AM, Alvarez Sedo CR, Blake D. Cleavage stage versus
573 blastocyst stage embryo transfer in assisted reproductive technology. *Cochrane Database Syst Rev*
574 2016:CD002118.
- 575 Harrell FE. *Regression Modeling Strategies: With Applications to Linear Models, Logistic Regression,*
576 *and Survival Analysis.* New York: Springer-Verlag , 2001.
- 577 Harrell FE, Lee KL, Mark DB. Multivariable prognostic models: issues in developing models, evaluating
578 assumptions and adequacy, and measuring and reducing errors. *Stat Med* 1996;**15**:361–387.
- 579 Leushuis E, van der Steeg JW, Steures P, Bossuyt PMM, Eijkemans MJC, van der Veen F, Mol BWJ,
580 Hompes PGA. Prediction models in reproductive medicine: a critical appraisal†. *Hum Reprod*

- 581 *Update* 2009;**15**:537–552.
- 582 Liber AC, Warner KE. Has Underreporting of Cigarette Consumption Changed Over Time? Estimates
583 Derived From US National Health Surveillance Systems Between 1965 and 2015. *Am J Epidemiol*
584 2018;**187**:113–119.
- 585 van Loendersloot L, Repping S, Bossuyt PMM, van der Veen F, van Wely M. Prediction models in in
586 vitro fertilization; where are we? A mini review. *J Adv Res* 2014;**5**:295–301.
- 587 van Loendersloot LL, van Wely M, Repping S, van der Veen F, Bossuyt PMM. Templeton prediction
588 model underestimates IVF success in an external validation. *Reprod Biomed Online* 2011;**22**:597–
589 602.
- 590 Malizia BA, Hacker MR, Penzias AS. Cumulative Live-Birth Rates after In Vitro Fertilization. *N Engl J*
591 *Med* 2009;**360**:236–243.
- 592 McLernon DJ, Steyerberg EW, te Velde ER, Lee AJ, Bhattacharya S. Predicting the chances of a live
593 birth after one or more complete cycles of in vitro fertilisation: population based study of linked
594 cycle data from 113 873 women. *BMJ* 2016;**355**:i5735.
- 595 Moons KGM, Kengne AP, Grobbee DE, Royston P, Vergouwe Y, Altman DG, Woodward M. Risk
596 prediction models: II. External validation, model updating, and impact assessment. *Heart*
597 2012;**98**:691–8.
- 598 Moons KGM, Kengne AP, Woodward M, Royston P, Vergouwe Y, Altman DG, Grobbee DE. Risk
599 prediction models: I. Development, internal validation, and assessing the incremental value of a new
600 (bio)marker. *Heart* 2012;**98**:683–90.
- 601 National Institute for Health and Care Excellence (NICE). Fertility problems: assessment and treatment.
602 Clinical guideline. 2013.

- 603 Nelson SM, Lawlor DA. Predicting live birth, preterm delivery, and low birth weight in infants born from
604 in vitro fertilisation: A prospective study of 144,018 treatment cycles. *PLoS Med* 2011;**8**:e1000386.
- 605 De Neubourg D, Bogaerts K, Blockeel C, Coetsier T, Delvigne A, Devreker F, Dubois M, Gillain N,
606 Gordts S, Wyns C. How do cumulative live birth rates and cumulative multiple live birth rates over
607 complete courses of assisted reproductive technology treatment per woman compare among
608 registries? *Hum Reprod* 2016;**31**:93–99.
- 609 Oakley L, Doyle P, Maconochie N. Lifetime prevalence of infertility and infertility treatment in the UK:
610 results from a population-based survey of reproduction. *Hum Reprod* 2008;**23**:447–450.
- 611 Olivius C, Friden B, Borg G, Bergh C. Why do couples discontinue in vitro fertilization treatment? A
612 cohort study. *Fertil Steril* 2004;**81**:258–61.
- 613 Oudshoorn SC, van Tilborg TC, Eijkemans MJC, Oosterhuis GJE, Friederich J, van Hooff MHA, van
614 Santbrink EJP, Brinkhuis EA, Smeenk JMJ, Kwee J, *et al*. Individualized versus standard FSH
615 dosing in women starting IVF/ICSI: an RCT. Part 2: The predicted hyper responder. *Hum Reprod*
616 2017;**32**:2506–2514.
- 617 Pandian Z, Marjoribanks J, Ozturk O, Serour G, Bhattacharya S. Number of embryos for transfer
618 following in vitro fertilisation or intra-cytoplasmic sperm injection. *Cochrane database Syst Rev*
619 2013;**7**:CD003416.
- 620 Rajkhowa M, McConnell A, Thomas GE. Reasons for discontinuation of IVF treatment: a questionnaire
621 study. *Hum Reprod* 2006;**21**:358–363.
- 622 Rossi B V, Berry KF, Hornstein MD, Cramer DW, Ehrlich S, Missmer SA. Effect of Alcohol
623 Consumption on In Vitro Fertilization. *Obstet Gynecol* 2011;**117**:136–142.
- 624 Rubin DB. Multiple Imputation for Nonresponse in Surveys. In: John Wiley & Sons, 2004.

- 625 Smeenk JM, Stolwijk AM, Kremer JA, Braat DD. External validation of the templeton model for
626 predicting success after IVF. *Hum Reprod* 2000;**15**:1065–8.
- 627 Smith ADAC, Tilling K, Lawlor DA, Nelson SM. External Validation and Calibration of IVFpredict: A
628 National Prospective Cohort Study of 130,960 In Vitro Fertilisation Cycles. Sun Q-Y (ed). *PLoS*
629 *One* 2015;**10**:e0121357.
- 630 Sterne JAC, White IR, Carlin JB, Spratt M, Royston P, Kenward MG, Wood AM, Carpenter JR. Multiple
631 imputation for missing data in epidemiological and clinical research: potential and pitfalls. *BMJ*
632 2009;**338**:b2393.
- 633 Steyerberg EW. *Clinical Prediction Models: A Practical Approach to Development, Validation, and*
634 *Updating*. New York, NY: Springer New York, 2009.
- 635 Stockwell T, Zhao J, Greenfield T, Li J, Livingston M, Meng Y. Estimating under- and over-reporting of
636 drinking in national surveys of alcohol consumption: identification of consistent biases across four
637 English-speaking countries. *Addiction* 2016;**111**:1203–1213.
- 638 Templeton A, Morris JK, Parslow W. Factors that affect outcome of in-vitro fertilisation treatment.
639 *Lancet* 1996;**348**:1402–1406.
- 640 van Tilborg TC, Eijkemans MJ, Laven JS, Koks CA, de Bruin JP, Scheffer GJ, van Golde RJ, Fleischer
641 K, Hoek A, Nap AW, *et al*. The OPTIMIST study: optimisation of cost effectiveness through
642 individualised FSH stimulation dosages for IVF treatment. A randomised controlled trial. *BMC*
643 *Womens Health* 2012;**12**:29.
- 644 van Tilborg TC, Oudshoorn SC, Eijkemans MJC, Mochtar MH, van Golde RJT, Hoek A, Kuchenbecker
645 WKH, Fleischer K, de Bruin JP, Groen H, *et al*. Individualized FSH dosing based on ovarian reserve
646 testing in women starting IVF/ICSI: a multicentre trial and cost-effectiveness analysis. *Hum Reprod*
647 2017;**32**:2485–2495. November 30, 2017.

- 648 van Tilborg TC, Torrance HL, Oudshoorn SC, Eijkemans MJC, Koks CAM, Verhoeve HR, Nap AW,
649 Scheffer GJ, Manger AP, Schoot BC, *et al.* Individualized versus standard FSH dosing in women
650 starting IVF/ICSI: an RCT. Part 1: The predicted poor responder. *Hum Reprod* 2017;**32**:2496–2505.
- 651 te Velde ER, Nieboer D, Lintsen AM, Braat DDM, Eijkemans MJC, Habbema JDF, Vergouwe Y.
652 Comparison of two models predicting IVF success; the effect of time trends on model performance.
653 *Hum Reprod* 2014;**29**:57–64.
- 654 Waylen AL, Metwally M, Jones GL, Wilkinson AJ, Ledger WL. Effects of cigarette smoking upon
655 clinical outcomes of assisted reproduction: a meta-analysis. *Hum Reprod Update* 2009;**15**:31–44.
- 656 Wong KM, Mastenbroek S, Repping S. Cryopreservation of human embryos and its contribution to
657 in vitro fertilization success rates. *Fertil Steril* 2014;**102**:19–26.
- 658 Zegers-Hochschild F, Adamson GD, Dyer S, Racowsky C, de Mouzon J, Sokol R, Rienzi L, Sunde A,
659 Schmidt L, Cooke ID, *et al.* The International Glossary on Infertility and Fertility Care, 2017†‡§.
660 *Hum Reprod* 2017;**32**:1786–1801.
- 661

662 **Tables**

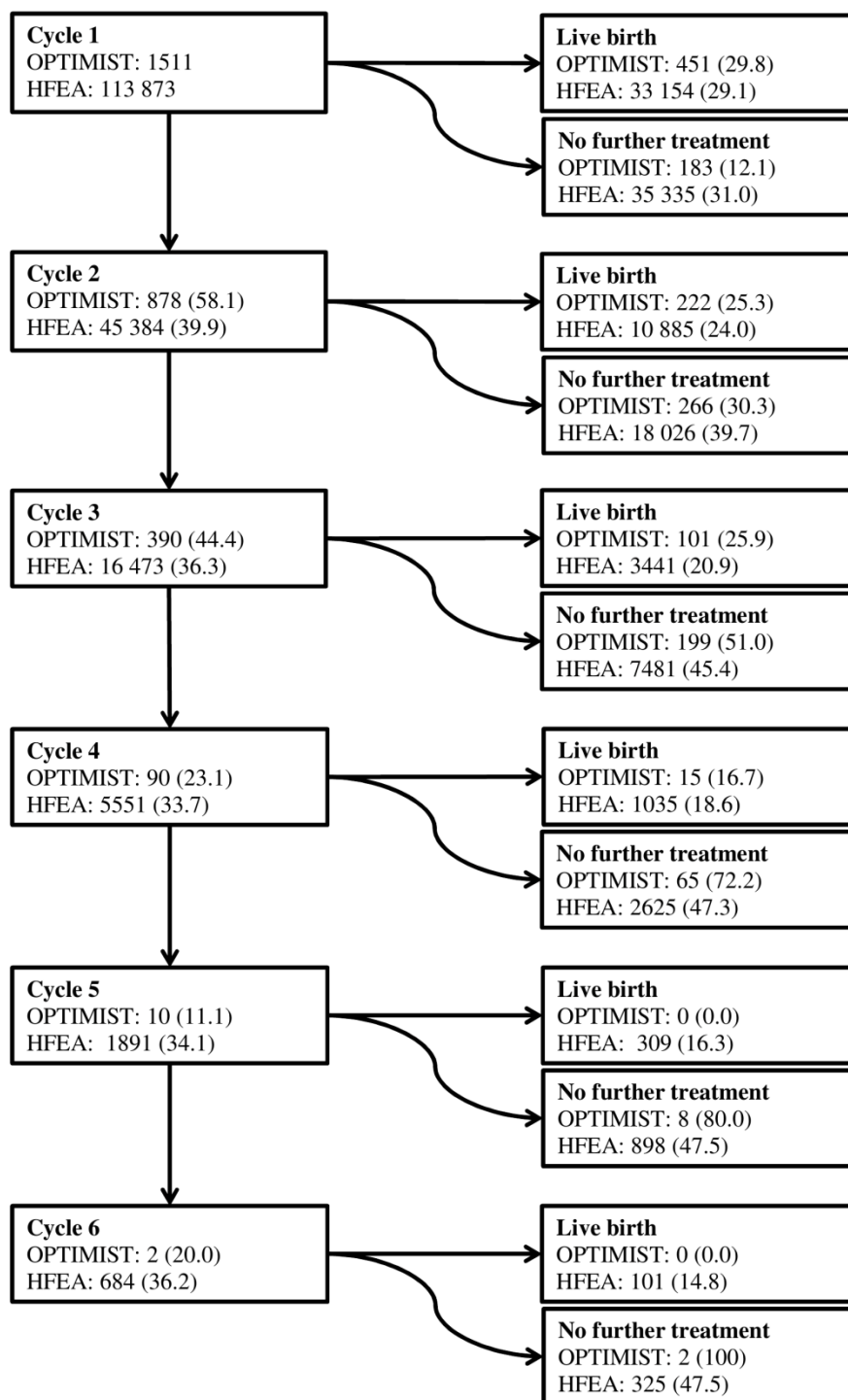
663 **Table I** Characteristics of patient and treatment variables included as predictors in the development
 664 sample (HFEA cohort) and the validation sample (OPTIMIST cohort) (McLernon *et al.*, 2016).

Characteristics	HFEA cohort	OPTIMIST cohort	Missing values in OPTIMIST cohort (%)
No of women	113 873	1 511	
No of complete cycles	184 269	2 881	
Patient characteristics			
Age (years), mean (SD)	34.1 (5)	33.5 (5)	2 (0.1)
Duration of infertility (years), median (IQR)	4 (3-6)	2 (2-3)	18 (1.2)
No previous pregnancy in couple	75 541 (66)	917 (61)	2 (0.1)
Cause of infertility:			
- Tubal factor	26 545 (23)	158 (11)	
- Male factor	49 753 (44)	839 (56)	
- Anovulatory	15 942 (14)	NA by protocol	
- Endometriosis	7 590 (7)	60 (4)	
- Unexplained	32 693 (29)	521 (35)	
Body weight (kg), mean (SD)	NA	69.5 (13)	36 (2.4)
Anti-Müllerian hormone (ng/mL), median (IQR)	NA	1.9 (1-3)	169 (11.2)
Antral follicle count (2-10mm), median (IQR)	NA	13 (9-18)	
Treatment characteristics of first completed cycle			
IVF	67 511 (59)	830 (55)	
ICSI	46 362 (41)	681 (45)	
No of oocytes retrieved, median (IQR)	8 (5-13)	8 (5-13) ^a	1 (0.1)
No of embryos created, median (IQR)	5 (2-8)	4 (2-7) ^a	4 (0.3)
No of embryos frozen, median (IQR)	0 (0-1)	1 (0-3) ^a	6 (0.5)
Cryopreservation of embryos	28 950 (25)	726 (48)	
Fresh embryo transfer: stage and no. of transferred embryos:			24 (1.6)
- Cleavage stage SET	9 248 (8)	1 004 (66)	
- Cleavage stage DET	75 701 (66)	125 (8)	
- Cleavage stage TET	8 649 (8)	4 (0.3)	
- Blastocyst stage SET	662 (1)	NA	
- Blastocyst stage DET	2 960 (3)	NA	
- Blastocyst stage TET	130 (0.1)	NA	
- No transfer	15 501 (14)	354 (23)	

665 Data are presented as number (%) unless otherwise specified. IQR; interquartile range, NA; not available, SET;
 666 single embryo transfer, DET; double embryo transfer, TET; triple embryo transfer.
 667 a) Median is calculated over 1293 women who had an ovarian follicle aspiration.

668 **Figures**

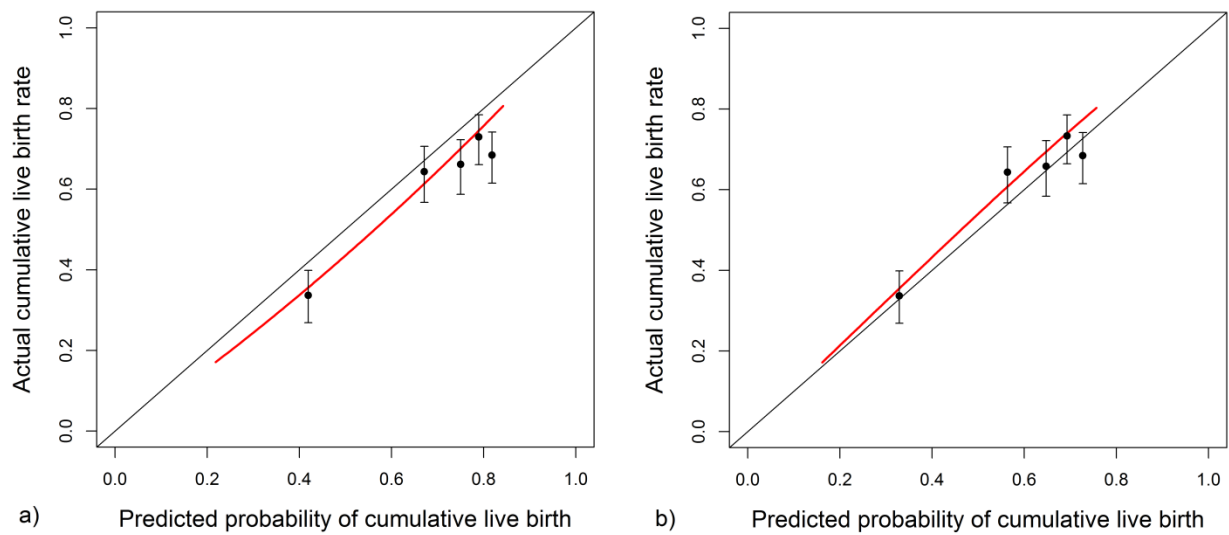
669 **Figure 1:** Flow chart presenting the numbers (%) of live birth, treatment continuation and discontinuation
 670 over six complete cycles in the OPTIMIST and HFEA databases (McLernon *et al.*, 2016).



671

672

673 **Figure 2:** Calibration plots showing the association between the calculated and observed cumulative live
674 birth rates over 3 complete IVF/ICSI cycles in the OPTIMIST cohort for **a)** the *original pre-treatment*
675 *model* as described by McLernon et al (McLernon *et al.*, 2016) **b)** *recalibrated pre-treatment model* with
676 adjustment of the intercept.

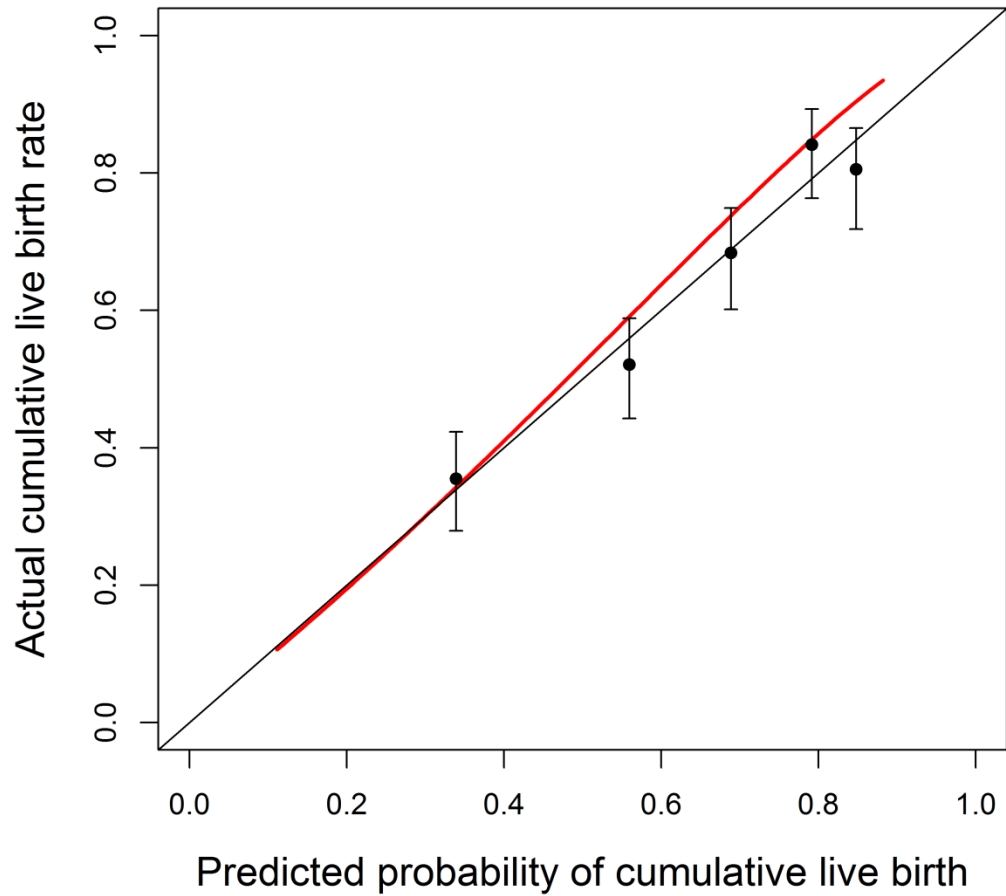


677

678

679

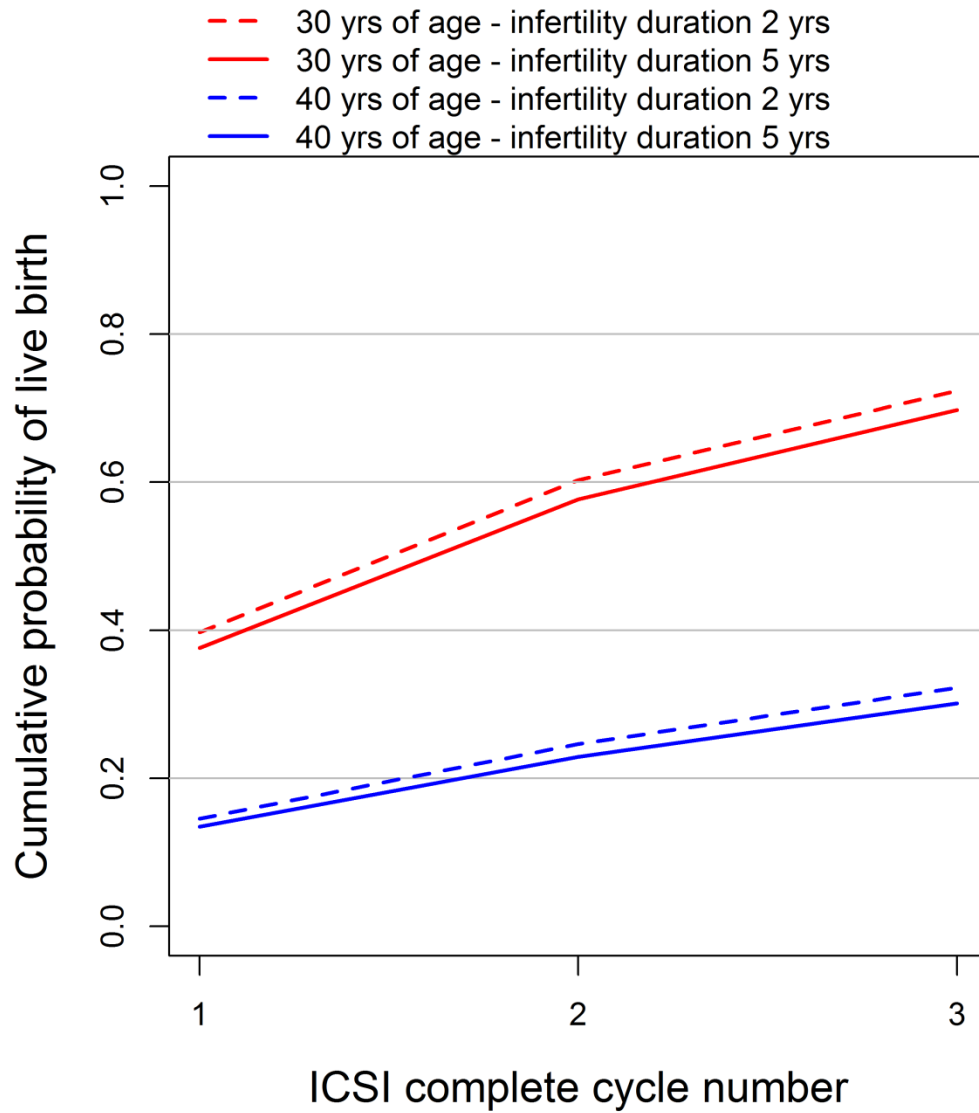
680 **Figure 3:** Calibration plot showing the association between the calculated and observed cumulative live
681 birth rates over 3 complete IVF/ICSI cycles in the OPTIMIST cohort for the *original post-treatment*
682 *model* as described by McLernon (McLernon *et al.*, 2016).



683

684

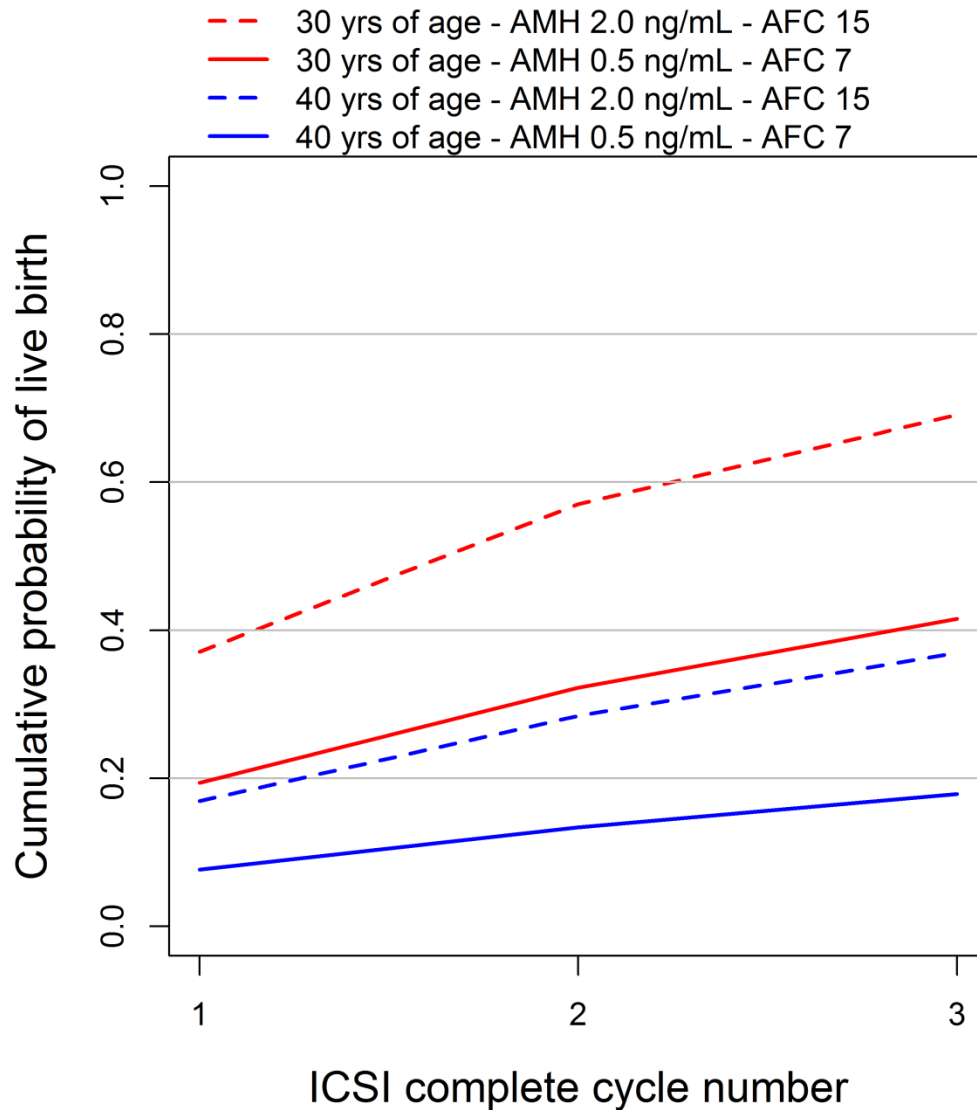
685 **Figure 4:** Example of the *recalibrated pre-treatment model* predicting the cumulative probability of a
686 live birth up to three complete ICSI cycles for a woman with primary infertility caused by a male factor,
687 aged 30 or 40 years with an infertility duration of two or five years.



688

689

690 **Figure 5:** Example of the with AMH, AFC and body weight *updated pre-treatment model* predicting the
 691 cumulative probability of a live birth up to three complete ICSI cycles for a woman with two years of
 692 primary infertility caused by a male factor, aged 30 or 40 years, a total body weight of 70 kilograms, with
 693 an AMH of 2.0 or 0.5 ng/mL and an AFC of 15 or 7.

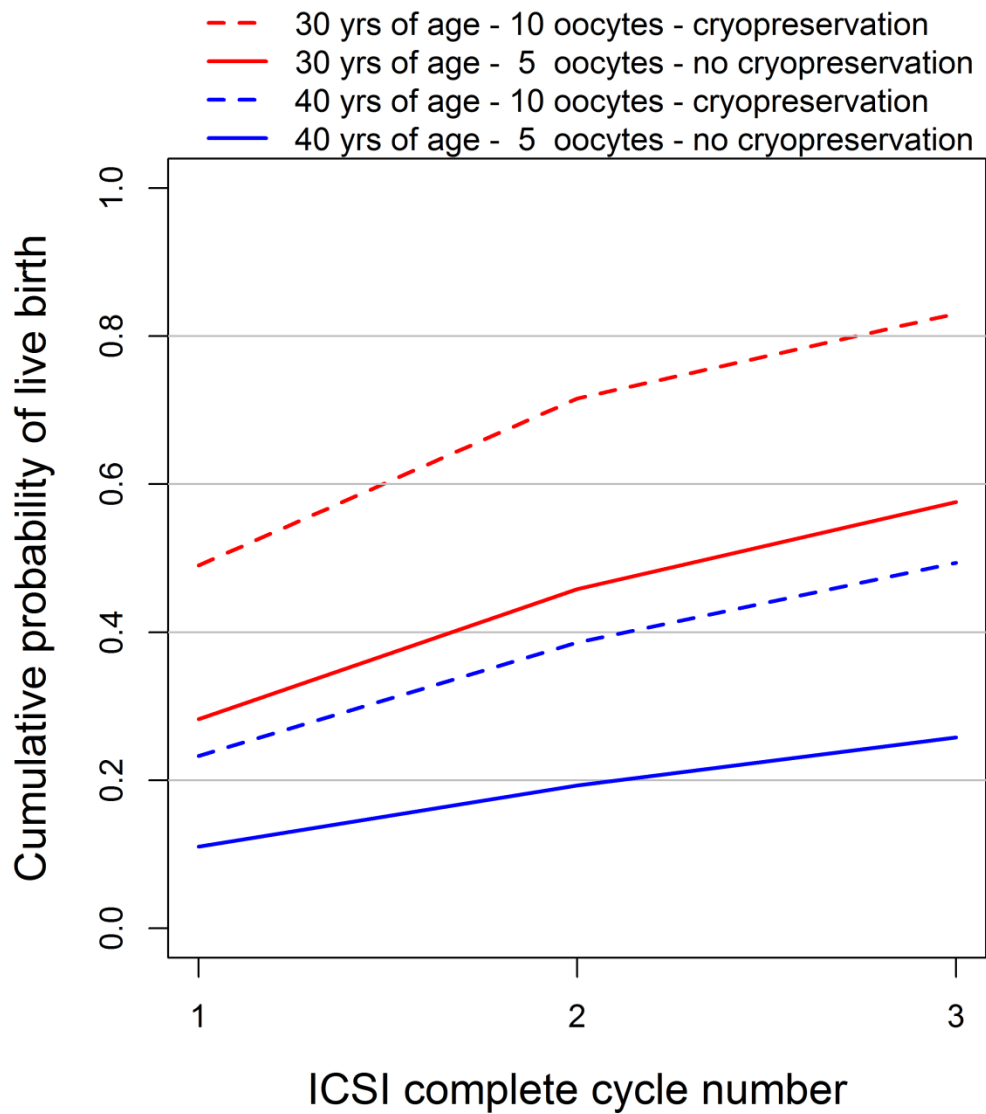


694

695

696 **Figure 6:** Example of the *post-treatment model* predicting the cumulative probability of a live birth up to
 697 three complete ICSI cycles for a woman with two years of primary infertility caused by a male factor,

698 aged 30 or 40 years, with 5 or 10 oocytes retrieved, a cleavage stage single embryo transfer, with or
699 without embryo cryopreservation.

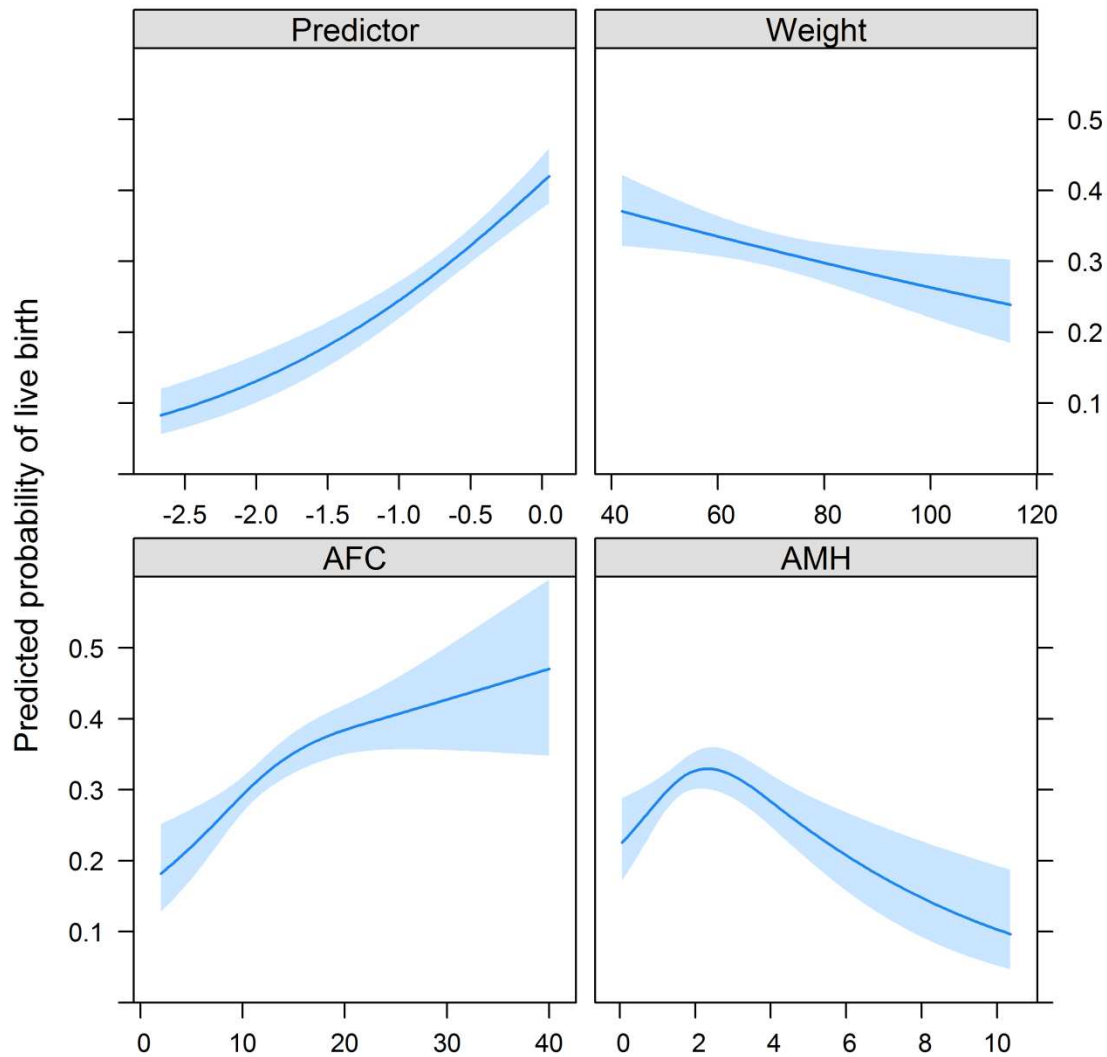


700

701

702 **Supplementary Figure 1.** Plots showing the adjusted relation between the predictors included in the
703 *updated McLernon pre-treatment model* and the probability of a live birth after IVF/ICSI treatment.

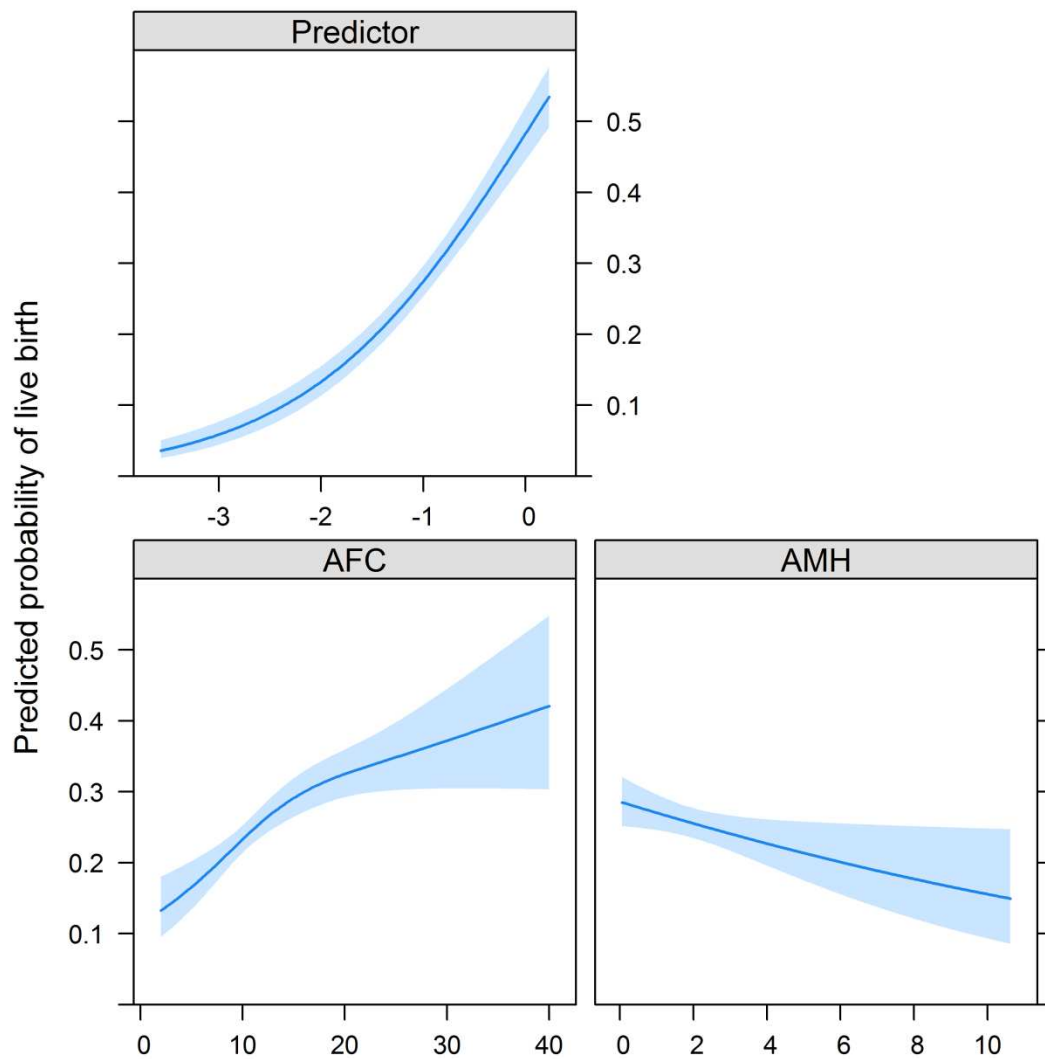
704 Predictor; linear predictor (XB) of the original pre-treatment model as described by McLernon
705 (McLernon et al. 2016), Weight; female body weight in kg, AFC; antral follicle count (2-10mm), AMH;
706 anti-Müllerian hormone (ng/mL)



707

708

709 **Supplementary Figure 2.** Plots showing the adjusted relation between the predictors in the *updated*
710 *McLernon post-treatment model* and the probability of a live birth after IVF/ICSI treatment.
711 Predictor: linear predictor (XB) of the original post-treatment model as described by McLernon
712 (McLernon et al 2016); AFC; antral follicle count (2-10mm), AMH; anti-Müllerian hormone (ng/mL)
713



714

715