# Journal Pre-proof

Predicting the outcomes of assisted reproductive technology treatments: A systematic review and quality assessment of prediction models

Ian Henderson, MSc, Michael P. Rimmer, MSc, Stephen D. Keay, FRCOG, Paul Sutcliffe, PhD, Khalid S. Khan, FRCOG, Ephia Yasmin, PhD, Bassel H. Al Wattar, PhD

Please cite this article as: Henderson I, Rimmer MP, Keay SD, Sutcliffe P, Khan KS, Yasmin E, Al Wattar BH, Predicting the outcomes of assisted reproductive technology treatments: A systematic review and quality assessment of prediction models, *F&S Reviews* (2021), doi: https://doi.org/10.1016/j.xfnr.2020.11.002.

1  **Predicting the outcomes of assisted reproductive technology treatments: A systematic**

2  **review and quality assessment of prediction models**

3

4  Ian Henderson MSc[1,2], Michael P Rimmer MSc[3], Stephen D Keay FRCOG[2], Paul Sutcliffe

5  PhD[1], Khalid S Khan FRCOG[4], Ephia Yasmin PhD[5], Bassel H.Al Wattar PhD[1,2,5]

6

7  [1]Warwick Medical School, Warwick University, Coventry, UK.

8  [2]Centre for Reproductive medicine, University Hospital Coventry and Warwickshire, Clifford

9  Bridge Road, Coventry, UK.

10  [3]MRC Centre for Reproductive Health, Queens Medical Research Institute, Edinburgh

11  BioQuarter, University of Edinburgh, UK.

12  [4]Department of Preventive Medicine and Public Health, University of Granada, 18071

13  Granada, Spain.

14  [5]Reproductive medicine unit, University College London Hospitals, London, UK.

15

16

17  **Corresponding author**: Bassel H.Al Wattar - Warwick Medical School, Warwick

18  University, Coventry, UK. Email: dr.basselwa@gmail.com.

19

20

21  **Short title:** Predicting assisted conception outcomes

22

23 **Capsule (30):**

24 We reviewed and evaluated 120 prediction models published over the last 24 years. We

25 identified twelve externally validated models that could be used to advise couples undergoing

26 fertility treatments.

27

28

29 **Abstract (250)**:

30 **Objective**: Predicting the outcomes of assisted reproductive technology (ART) treatments is

31 desirable, but adopting prediction models into clinical practice remains limited. We aimed to

32 review available prediction models for ART treatments by conducting a systematic review of

33 the literature to identify the best performing models for their accuracy, generalisability and

34 applicability.

35 **Evidence review:** We searched electronic databases (MEDLINE, EMBASE, and

36 CENTRAL) until June 2020. We included studies reporting on the development or evaluation

37 of models predicting the reproductive outcomes before (pre-ART) or after starting (Intra-

38 ART) treatment in couples undergoing any ART treatment. We evaluated the models'

39 discrimination, calibration, type of validation, and any implementation tools for clinical

40 practice.

41 **Results**: We included 69 cohort studies reporting on 120 unique prediction models. Half the

42 studies reported on pre-ART (48%) and half on intra-ART (56%) prediction models. The

43 commonest predictors used were maternal age (90%), tubal factor subfertility (50%), and

44 embryo quality (60%).

45 Only fourteen models were externally-validated (14/120, 12%) including eight pre-ART

46 models (Templeton, Nelson, LaMarca, McLernon, Arvis, and the Stolwijk A/I,C,II models),

47 and five intra-ART models (Cai, Hunault, van Loendersloot, Meijerink, Stolwijk B, and the

48 McLernon post-treatment model) with a reported c-statistics ranging from 0.50 to 0.78. Ten

49 of these models provided implementation tools for clinical practice with only two reported

50 online calculators.

51 **Conclusion**: We identified externally validated prediction models that could be used to

52 advise couples undergoing ART treatments on their reproductive outcomes. The quality of

3

53  available models remains limited and more research is needed to improve their

54  generalizability and applicability into clinical practice.

55
56  **Keywords:** infertility, prediction, assisted reproduction, systematic review.

57

58  **Highlights:**

59  -   Over the last 24 years a high number of studies attempted to produce useful prediction

60      models and decision aids for clinicians and patients undergoing ART.

61  -   In this review we evaluated 69 studies reporting on 120 unique prediction models, but

62      only a minority of these models were externally validated or useful in clinical

63      practice.

64  -   Most of these models suffered from a high risk of bias driven by poor model

65      development, data sampling and analysis methodology.

66  -   More research is needed to leverage available data, refine published models, and

67      increase their applicability in clinical practice using novel technology such as

68      artificial intelligence and dynamic intra-treatment prediction modelling.

69

70

71 **Introduction**

72 Assisted reproductive technology (ART) has evolved over the last 40 years offering hope to a

73 record number of  infertile couples worldwide (1–3). Currently ART is the first port of call

74 for many couples inclusive of those experiencing unexplained and reversible causes of

75 subfertility such as mild male factor and unilateral tubal pathology. The birth rate with

76 assisted conception increased steadily over the last few decades from an average of 9% in

77 1991 to 23% in 2018 (4). This mass adoption of ART, however, sparked the debate on the

78 ethical use of some ART treatments (5), their cost-effectiveness, and the risk of profiteering

79 to certain patient groups (6). Accurate prediction of clinical outcomes and any mitigating risk

80 factors could help to rationalize the use of ART treatments and improve their clinical

81 effectiveness (7). While many prediction models have been produced to aid clinicians and

82 couples in planning their fertility treatments, implementing those models remains limited in

83 practice (8).

84

85 To be used effectively, prediction models should undergo rigorous development, validation,

86 and impact assessment (9,10). Unsurprisingly, few published models complete this process

87 which limits their clinical value and increase research wastage (7,8,11). Advances in data

88 gathering and statistical methodology using machine learning and artificial intelligence could

89 help to streamline the development and validation process of prediction models, but such

90 practice remains limited in reproductive medicine (12).

91

92 Our aim was to systematically review and evaluate the performance, generalisability and

93 applicability of published prediction models for ART treatments to identify the best

94 performing models that could be used in clinical practice.

95

5

96    **Methods**

97    We conducted this systematic review using a prospectively registered protocol

98    (CRD42019156606) and reported the findings following standard guidelines (13).

99

100    *Search strategy and study selection*

101    We searched electronic databases (MEDLINE, EMBASE, and Cochrane CENTRAL) from

102    inception until June 2020 for all studies reporting on the development or evaluation of any

103    prediction model for the outcome of any ART treatments (in vitro fertilization (IVF) and/or

104    intracytoplasmic sperm injection (ICSI)). We did not apply any search filters or language

105    restrictions. Articles in non-English were translated if deemed relevant. We conducted

106    supplementary searches in Google Scholar and Scopus for any additional articles of interest

107    in the grey literature. We also searched the bibliographies of relevant articles to identify any

108    missing citations.

109

110    We included longitudinal studies that reported on the development or evaluation of any

111    model for predicting clinical pregnancy (confirmed on ultrasound) or live birth following any

112    ART treatments. We excluded studies reporting on the crude association between a single

113    independent variable and the outcomes of interest, those reporting on non-predictive models,

114    and those not reporting on the model performance measures. Models predicting non-

115    reproductive outcomes or solely predicting biochemical pregnancy were also excluded.

116    Similarly, we excluded models that used solely embryological or seminal parameters to

117    predict the outcomes of interest. Finally, we also excluded case series, conference abstracts

118    and review articles.

119

120    *Assessment of study quality*

121 We assessed the risk of bias and applicability of the included studies in duplicate using the

122 PROBAST tool (14). Studies were assessed in four domains: population, predictors, outcome,

123 and analysis. Studies were deemed low risk of bias if they were cohort studies, defined and

124 measured predictors consistently and independently of the pre-specified outcome, included

125 sufficient events per variable with appropriate parameterisation of predictors, included all

126 participants in the analysis, treated missing data appropriately, did not include predictors

127 based on univariable analyses, assessed the model's discrimination and calibration

128 appropriately, and accounted for model overfitting and optimism based on the use of an

129 appropriate validation procedure and shrinkage of estimates in the presence of optimism

130 which were evaluated in the context of events per variable, appropriate parameterisation and

131 modelling strategy (14). We produced an overall assessment of both the risk of bias and

132 model applicability per study.

133

134 *Models performance, generalizability and applicability*

135 We evaluated models' performance by their reported discrimination (the model's ability to

136 separate those with and without the outcome of interest) and calibration (the concordance

137 between predicted and observed outcome frequency) measures (15). Discrimination is

138 commonly described using the rank order statistic 'area under the receiver operating

139 characteristic curve' (AUROC), which is equivalent to the concordance-statistic (c-statistic).

140 We considered a c-statistic value of 0.5 to represent no discriminative ability, a value of 1 to

141 represent perfect discriminative ability (15). Calibration is often assessed using the Hosmer-

142 Lemeshow statistic (16). A model is considered well-calibrated when the average predicted

143 probability per sub-group matches the observed proportion. Calibration is more informatively

144 assessed graphically by the calibration plot, where the predicted probability per ordered sub-

145 group is plotted against the observed proportion, demonstrating the nature and magnitude of

7

146    any miscalibration. An intercept of 0 and a slope of 1 therefore represents perfect calibration

147    (17).

148

149    To evaluate generalizability, we reported on the validation process for each model including

150    the validation type, procedures, and characteristics of the validation population. We divided

151    validation efforts into 'internal', 'temporal', or 'external' depending the type of validation

152    population.

153

154    To evaluate the models' applicability and translation into clinical practice, we reported on

155    efforts to increase the model's accessibility to both health professionals and lay consumers,

156    and the availability of any decision support tools including predicted probabilities based on

157    patient profile, score-based decision aids, score-based nomograms, to end-user web-based

158    predictive calculators.

159

160    *Data extraction*

161    Two independent reviewers (IH and MPR) extracted data onto a custom designed collection

162    database guided by the CHARMS checklist (18) to identify relevant data points for extraction

163    and reporting. We extracted data on the study design, outcome, sample size, population

164    characteristics, model development methods, performance and validation statistics, and

165    clinical application. We divided models into (pre-ART) where outcome prediction was

166    possible prior to commencing ovarian stimulation, and (intra-ART) where outcome

167    prediction was possible after commencing ovarian stimulation. We categorized the included

168    studies as per the TRIPOD guidelines into: type 1a studies developing a model and evaluating

169    its predictive performance using the same data (apparent performance), type 1b studies

170    developing a prediction model using the entire dataset with resampling (e.g. bootstrapping or

8

171 cross-validation) techniques to evaluate the performance and optimise the developed model,

172 type 2a studies with data randomly split to develop the model and then to evaluate its

173 predictive performance, type 2b studies with data non-randomly split (e.g. by location or

174 time) to develop the prediction model and then to evaluate its predictive performance, type 3

175 studies developing a prediction model using one dataset and an evaluation of its performance

176 on separate data (e.g. from a different population), and type 4 studies which are only

177 evaluating the predictive performance of an existing prediction model in a separate dataset

178 (19).

179

180 *Statistical analysis*

181 We summarised data using descriptive statistics and reported on continuous data using means

182 or medians with standard deviations where relevant. For dichotomous data we reported using

183 frequencies and natural percentages. All analyses and figures were produced using RStudio

184 version 1.2.1335 (RStudio, Boston, MA) (20).

185

186 **Results**

187 *Study characteristics*

188 Our search revealed 8052 potentially relevant unique citations; of these, we reviewed 483 in

189 full and included 69 studies in our review reporting on the development of 120 ART

190 prediction models (Figure 1). All included studies were cohort studies, 55 of which were

191 retrospective (55/69, 79.7%) and 14 prospective (14/69, 20.3%). As per TRIPOD

192 classification, 18 (18/69, 26.1%) of these studies were type 1a studies, 20 (20/69, 29.0%)

193 were type 1b, 6 (6/69, 8.7%) were type 2a, 10 (10/69, 14.5%) were type 2b, 5 (5/69, 7.2%)

194 were type 3, and 10 (10/69, 14.5%) type 4 (Figure 2). The majority were from Europe (49/69,

195 71.0%) with only eleven from Asia (11/69, 15.9%), and three from North America (3/69,

196 4.3%).

197

198 There were variations in the population characteristics across included studies. Nine studies

199 (13.4%) included unselected couples (for age, cycle cancellation, maternal comorbidity,

200 aetiology, and sperm source), seven included unselected couples but excluded women using

201 donor gametes (10.4%), and twelve studies (17.9%) included couples with selected baseline

202 characteristics (Supplementary Table 1). About half of the included studies explicitly

203 excluded donor oocyte cycles (29/69, 42.0%), and a third explicitly excluded cancelled cycles

204 (21/69, 30.4%), and a quarter explicitly excluded women outside a specific age range (18/69,

205 26.1%).

206

207 Most of the included studies reported on the development (with or without validation) of

208 novel models (62/69, 89.9%), with the remainder uniquely reporting on the validation of pre-

209 existing models (7/69, 10.1%). Half of these studies (30/62, 48.3%) reported on pre-ART

210 predictive models (21–47), and 56% (35/62, 56.5%) reported on intra-ART (48–78). Only

211 three studies (3/62, 4.8%) reported on both pre and intra-ART predictive models (79–81).

212 Three quarters of these developmental studies (47/62, 75.8%) involved IVF/ICSI treatments,

213 twelve IVF treatment only (12/62, 19.4%), and two ICSI treatment only (2/62, 3.2%), with 1

214 unspecified by the authors. Two-thirds included only cycles using a fresh embryo transfer

215 (41/62, 66.1%), while both fresh and frozen embryo cycles were included in 21 studies

216 (21/62, 33.9%).

217

218 *Predictors and outcomes*

10

219   For studies that developed pre-ART models, the commonest included predictor was maternal

220   age (27/30, 90.0%) followed by tubal factor subfertility (15/30, 50.0%), gravidity (13/30,

221   43.3%), and the duration of subfertility (12/30, 40.0%) (Figure 3a). A similar trend was seen

222   for intra-ART models as the commonest included predictor was also maternal age (33/35,

223   94.3%), followed by embryo quality (21/35, 60.0%), previous ART success (16/35, 45.7%),

224   duration of subfertility (12/35, 34.3%), and tubal factor subfertility (10/35, 28.6%) (Figure

225   3b).

226

227   Live birth was the outcome of interest across all studies, for those that developed both pre-

228   ART (20/30, 66.7%) and intra-ART (18/35, 51.4%) models. A quarter of studies that

229   developed intra-ART models focused on clinical pregnancy (10/35, 28.6%) and ongoing

230   pregnancy (8/35, 22.9%) which were less frequently reported in pre-ART models (clinical

231   pregnancy (5/30, 16.7%), ongoing pregnancy (5/30, 16.7%)).

232

233   *Sample size and modelling method*

234   The median sample size for developing pre-ART models was 757 for participants (range 85-

235   113,873) and 1,061 for ART cycles (range 113-443,202). For intra-ART models, the median

236   participant sample size was 1,419 (range 90-113,873) and median ART cycles was 1,676

237   (range 110-184,269). Most studies (48/69, 69.6%) had ≥10 events per candidate variable

238   (degrees of freedom). The majority of studies developed models using logistic regression

239   (pre-ART (24/30, 80.0%), intra-ART (30/35, 85.7%)). Only a minority used other methods,

240   including generalized estimating equations, Bayesian networks, Cox regression, machine

241   learning techniques and deep learning techniques (Supplementary Table 2).

242

243   *Performance, generalizability and applicability*

11

244   Discrimination was reported for most of the included studies (109/120, 90.8%) while

245   calibration was reported for over half (72/120, 60.0%). Both discrimination and calibration

246   were reported in only 61 studies (61/120, 50.8%). The commonest methods to assess

247   calibration were the Hosmer-Lemeshow statistic (27/72, 37.5%), calibration plot (24/72,

248   33.3%), slope test (14/72, 19.4%), and calibration-in-the-large (11/72, 15.3%).

249

250   We captured 31 unvalidated models from type 1a studies without subsequent validation

251   (31/120, 25.8%), as well as six models that were locally refit from validation studies (6/120,

252   5.0%). Fifty-five models were internally-validated from 1b/2a studies without subsequent

253   validation (55/120, 45.8%), 15 were temporally-validated models from 2b studies without

254   subsequent validation (15/120, 12.5%). There were seven external validation studies (7/120,

255   5.8%). Four were type 4 studies by a team that overlapped with the model development team

256   (4/120, 3.3%)(35,80,82),(22,23,79), and three studies were performed by independent

257   validation teams (30,37,57)

258   We captured eight externally validated pre-ART models: the Templeton model (n=6

259   validations), Nelson model (n=3), LaMarca model (n=1), McLernon pre-treatment model

260   (n=1), Arvis model (n=1), and The Stolwijk models A/I, C, and II (n=7). All models showed

261   similar performance with c-statistics ranging from 0.53 to 0.78. The Stolwijk models A/I and

262   II were declared invalid (Table 1).

263

264   Among the intra-ART models, only five were externally validated: the Cai model (n=1),

265   Hunault model (n=1), van Loendersloot model (n=1), Meijerink model (n=1), and the

266   McLernon post-treatment model (n=1). All models showed similar performance with c-

267   statistics ranging from 0.63 to 0.78. However, only the McLernon model was validated in a

268 good quality external validation study with low risk of bias showing a c-statistic of 0.71

269 (95%CI 0.69-0.74) and reportedly good calibration (Table 1).

270

271 Only a quarter of all published models (33/120, 25.4%) were presented in full either offering

272 the regression formula, coefficients with intercept, or baseline hazard. Seven models

273 presented nomograms or score charts (7/120, 5.8%), and seven were adapted into online risk

274 prediction calculators (7/120, 5.8%). Of these, only three calculators were functional at the

275 time of writing this review(83–85). Overall, half of the included studies (35/62, 56.5%),

276 reporting on 47 models (47/120, 39.2%), enabled the reader to generate a personalised

277 prediction in a useful format. All the externally validated models offered an implementation

278 tool except the Cai model and the invalid Stolwijk models. But only two presented an online

279 calculator for use by health professionals and patients (the Nelson and the McLernon

280 calculators) (Table 1).

281

282 *Quality and risk of bias*

283 Overall, a majority of the included studies were at high risk of bias (56/69, 81.2%) and only

284 ten studies at low risk (10/69, 14.5%) (Figure 4, Supplementary Table 3). Within the

285 'participant' domain, three-quarters of the included studies were at low risk (50/69, 72.5%)

286 and nine at high risk (9/69, 13.0%). Similarly, within the 'outcome' domain, the majority

287 were at low risk (66/69, 95.7%). In contrast, within the 'predictor' domain only half were at

288 low risk (32/69, 46.4%), with 36 studies of unclear risk due to providing inadequate

289 definitions, namely for candidate predictors (36/69, 52.2%). For the 'analysis' domain, less

290 than a fifth were of low risk of bias (12/69, 17.4%). Half (35/69, 50.7%) assessed model

291 performance appropriately, by discrimination and an informative measure of calibration.

292 Only a quarter reported and handled missing data appropriately (16/69, 23.2%); only 19

293  studies (19/69, 27.5%) addressed overfitting and optimism; only 48 had sufficient events per

294  candidate predictor (≥20 events (14)) (48/69, 69.6%), and only 38 parameterized predictors

295  appropriately (38/69, 55.1%).

296

297  **Discussion**

298  *Summary of main findings*

299  Our findings depict an overall high investment in producing working prediction models and

300  decision aids for clinicians and patients undergoing ART treatments with 120 models

301  produced over the last 24 years, an average of 5 models produced per year. However, while

302  huge resources and patient data were committed to producing these models, only a minority

303  of these studies offered externally validated models that could be used in everyday practice.

304

305  The majority of the included studies had a high risk of bias, largely driven by poor model

306  development methodology specifically in data sampling and analysis (Figure 4). Only a

307  minority of models were developed within large sizes cohorts (only 9 studies included

308  >10,000 women/cycles) and most were selected ART populations, thus reducing model's

309  applicability in practice. In contrast, with much prediction data available several clinical and

310  biochemical markers are now well established as reliable predictors of reproductive outcomes

311  (Figure 3a, 3b). Leveraging this large body of evidence could facilitate the process of

312  developing and validating future models to minimize duplication of efforts. Logistic

313  regression modelling remains the commonest method for model development, though

314  alternative methodology is becoming popular such as artificial intelligence aided techniques

315  (29,34,38,46,48,49,54,65,69,75,86).

316

317  *Strengths and limitations*

318  The strengths of our review are several. In contrast to previously published reviews (7,8,11),

319  we used a prospectively registered protocol, applied  a comprehensive search strategy,

320  extracted data in duplicate, assessed quality according to PROBAST criteria, and included all

321  types of studies as per TRIPOD (both model development and validation studies) to evaluate

322  models' applicability into clinical practice. Consequently, our findings offer a robust

323  assessment of the current state-of-the-art in ART prediction modelling and the remaining

324  knowledge gap. To aid their adoption in practice, we identified top performing models

325  referencing their quantitative assessment markers, relevant population of interest and how

326  they can be accessed online (Table 1).

327

328  Our research was inclusive with almost double the number of studies included in the most

329  recent review (11) offering a more comprehensive and systematic assessment of the

330  literature. A previous review by Ratna et al adopted an arbitrary quality threshold of 80%

331  adherence to TRIPOD (19) in their inclusion criteria which could have limited the

332  generalizability of their findings. We refrained from imposing any reporting thresholds and

333  assessed the methodological quality of all published models to offer a comprehensive and

334  objective assessment of the literature.

335

336  Our findings still have some limitations. Several of the studies reported vaguely on the

337  measures of calibration using terms like "good calibration" which limited our ability to

338  provide an objective assessment of these models. Furthermore, given the lack of a universally

339  adopted definition of what constitutes good calibration for ART models, it is difficult to

340  preferentially select top performing models. Clearly, most subfertile couples have some

341  probability of conceiving independent of any treatment, similarly the chance of conception in

342  healthy couples is never 100% in every cycle. As the methodological standards for model

15

343 development improved over time, our contemporary PROBAST assessment of risk of bias

344 might differ from older reviews and the findings are therefore not completely reproducible.

345

346 *Implications for clinical practice*

347 Introducing prediction modelling into clinical practice was aimed to tailor treatments to each

348 patient's individual needs, thus maximising effectiveness and reducing personal harm (9).

349 Models can aid decision making on starting treatment (87) or to adjust a treatment to the

350 patient characteristics (88). Whilst most treatments are static (e.g., medication or surgery), the

351 process of undergoing IVF or ICSI treatments is heterogeneous and dynamic, continuously

352 changing through a series of interconnected complex decisions made to optimise successful

353 conception. Coupled with the rapid progress in ART, it is likely that most models will be

354 over-simplistic and become outdated. This applies especially to pre-ART models which are

355 dependent on a limited range of predictors that cannot adjust for initial treatment response

356 (e.g., ovulation stimulation and embryo fertilisation). Consequently, the clinical value of

357 available models is currently limited to counselling patients on the value of starting ART

358 treatment rather than tailoring those treatments to maximize chances of conception. A

359 solution could lie in the development, validation and continuous update of dynamic models

360 that could adjust for the within-treatment changes and offer a refined estimate of successful

361 conception throughout the ART treatment process (89).

362

363 The process of IVF/ICSI is emotionally and psychologically demanding with patients often

364 having to make difficult decisions such as the use of frozen embryos or consider add-on

365 therapies (90). Predicting the chances of conception in itself can be stressful (91) which could

366 limit the adoption of these models in practice. As such, developing any prediction models

367 should be guided by expressed patients' needs (92), a practice we did not observe in the

16

368  models included in this review. Future model development should take into account the

369  various decision-making processes involved in the ART treatment process and the associated

370  predictors that could add cumulative information to aid patients and their caring clinicians in

371  the decision-making process. Lastly, successful model implementation into clinical practice

372  could be facilitated by improved interpretability (93) and user-friendly interfaces that enable

373  end users to input and access data effortlessly in jargon-free outputs such as online risk

374  calculators or decision aid tools hosted on mobile apps (83–85).

375

376  *Future research need*

377  Our findings illustrate an abundance of data dedicated to predict ART outcomes, yet

378  translation into practice remains limited. As our ability to collect and analysis large datasets

379  improves over time, perhaps future steps should focus more on harmonizing data collection

380  across institutions, regulators and countries to facilitate streamlined model development,

381  validation, and update while reducing associated costs. Crucially, there is a need to focus

382  available resources on combining data from published models (e.g., using individual patient

383  data meta-analysis methodology) and externally validating ensuing ones rather than on

384  developing newer models.

385

386  We captured a recent trend towards using artificial intelligence (AI) technology in model

387  development (29,34,38,46,48,49,54,65,69,75,86). While promising, most of these models did

388  not achieve improved prediction performance nor followed sound methodology compared to

389  older ones (94). Specifically, the work on many of these models seem to be driven by an

390  experimental approach evaluating the different AI technologies rather than a multi-

391  disciplinary approach aiming to address real patients' needs. Still, leveraging the power of AI

392  technology and big data research methods to simulate the complex decision making process

393  involved in ART treatments could be a game changer to provide accurate individualized

394  fertility assessment to couples in need (95). Large multi-national multi-disciplinary teams are

395  best equipped to address this complex and important health problem.

396

397  **Conclusions**

398  We identified externally validated prediction models that could be used to advise couples

399  undergoing ART treatments on their reproductive outcomes. The quality of available models

400  remains limited and more research is needed to improve their generalisability and

401  applicability in clinical practice.

402

403

412

413  **Contribution to Authorship:** BHA conceived the idea. BHA and IH wrote the final protocol

414  and manuscript. IH conducted the search. IH and MR conducted the data extraction and 1[st]

415  draft of the manuscript. BHA and IH conducted the statistical analysis and data interpretation.

416  SK and KSK contributed to data interpretation and final editing of the manuscript.

417

418

419

420

421

422

423

424

**References**:

1. wwwhfeagovuk. HFEA Fertility treatment 2017: trends and figures [Internet]. 2017. Available from: www.hfea.gov.uk

2. Society for Assisted Reproductive Technology. National Summary Report [Internet]. 2017;Available from: https://www.sartcorsonline.com/rptCSR_PublicMultYear.aspx

3. Zegers-Hochschild F, Schwarze JE, Crosby J, Musri C, Urbina MT. Assisted reproductive techniques in Latin America: the Latin American Registry, 2015. Reprod Biomed Online 2018;37:685–92.

4. Human Fertilisation & Embryology Authority. Fertility treatment 2018: trends and figures. 2020.

5. te Velde E, Habbema D, Nieschlag E, Sobotka T, Burdorf A. Ever growing demand for in vitro fertilization despite stable biological fertility—A European paradox. Eur. J. Obstet. Gynecol. Reprod. Biol. 2017;214:204–8.

6. Heng BC. Can the difference in medical fees for self and donor freeze-thaw embryo transfer cycle, be in fact a cover-up for the sale of donated human embryos? Philos. Ethics, Humanit. Med. 2007;

7. Leushuis E, van der Steeg JW, Steures P, Bossuyt PMM, Eijkemans MJC, van der Veen F, et al. Prediction models in reproductive medicine: a critical appraisal. Hum Reprod Update 15:537–52.

8. van Loendersloot L, Repping S, Bossuyt PMM, van der Veen F, van Wely M. Prediction models in in vitro fertilization; where are we? A mini review. J. Adv. Res. 2014;

9. Steyerberg EW, Vickers AJ, Cook NR, Gerds T, Gonen M, Obuchowski N, et al. Assessing the performance of prediction models: A framework for traditional and novel measures. Epidemiology. 2010;21:128–38.

450    10.    Steyerberg EW, Harrell FE. Prediction models need appropriate internal, internal-

451         external, and external validation. J. Clin. Epidemiol. 2016;

452    11.    Ratna MB, Bhattacharya S, Abdulrahim B, McLernon DJ. A systematic review of the

453         quality of clinical prediction models in in vitro fertilisation. Hum Reprod

454         2020;35:100–16.

455    12.    Kelly CJ, Karthikesalingam A, Suleyman M, Corrado G, King D. Key challenges for

456         delivering clinical impact with artificial intelligence. BMC Med. 2019;17:1–9.

457    13.    Moher D, Liberati A, Tetzlaff J, Altman DG. Preferred reporting items for systematic

458         reviews and meta-analyses: The PRISMA statement. BMJ. 2009;339:332–6.

459    14.    Wolff RF, Moons KGM, Riley RD, Whiting PF, Westwood M, Collins GS, et al.

460         PROBAST: A tool to assess the risk of bias and applicability of prediction model

461         studies. Ann Intern Med 2019;170:51–8.

462    15.    Cook NR. Use and misuse of the receiver operating characteristic curve in risk

463         prediction. Circulation 2007;115:928–35.

464    16.    HOSMER DW, HOSMER T, CESSIE S LE, LEMESHOW S. A COMPARISON OF

465         GOODNESS□OF□FIT TESTS FOR THE LOGISTIC REGRESSION MODEL. Stat

466         Med 1997;16:965–80.

467    17.    Miller ME, Langefeld CD, Tierney WM, Hui SL, Mcdonald CJ. Validation of

468         Probabilistic Predictions. Med Decis Mak 1993;

469    18.    Moons KGM, de Groot JAH, Bouwmeester W, Vergouwe Y, Mallett S, Altman DG, et

470         al. Critical Appraisal and Data Extraction for Systematic Reviews of Prediction

471         Modelling Studies: The CHARMS Checklist. PLoS Med 2014;11:e1001744.

472    19.    Moons KGM, Altman DG, Reitsma JB, Ioannidis JPA, Macaskill P, Steyerberg EW, et

473         al. Transparent reporting of a multivariable prediction model for individual prognosis

474         or diagnosis (TRIPOD): Explanation and elaboration. Ann Intern Med 2015;162:W1–

21

475      73.

476   20.   RStudio Team. RStudio: Integrated Development for R. 2018;

477   21.   Alebić MŠ, Stojanović N, Zuvić-Butorac M. The IVF Outcome Counseling Based on

478      the Model Combining DHEAS and Age in Patients with Low AMH Prior to the First

479      Cycle of GnRH Antagonist Protocol of Ovarian Stimulation. Int J Endocrinol

480      2013;2013:637919.

481   22.   Arvis P, Lehert P, Guivarc'h-Levêque A. Simple adaptations to the Templeton model

482      for IVF outcome prediction make it current and clinically useful. Hum Reprod

483      2012;27:2971–8.

484   23.   La Marca A, Nelson SM, Sighinolfi G, Manno M, Baraldi E, Roli L, et al. Anti-

485      Müllerian hormone-based prediction model for a live birth in assisted reproduction.

486      Reprod Biomed Online 2011;22:341–9.

487   24.   Li HWR, Lee VCY, Lau EYL, Yeung WSB, Ho PC, Ng EHY. Role of Baseline Antral

488      Follicle Count and Anti-Mullerian Hormone in Prediction of Cumulative Live Birth in

489      the First In Vitro Fertilisation Cycle: A Retrospective Cohort Analysis. PLoS One

490      2013;8:e61095.

491   25.   Lintsen AME, Eijkemans MJC, Hunault CC, Bouwmans CAM, Hakkaart L, Habbema

492      JDF, et al. Predicting ongoing pregnancy chances after IVF and ICSI: a national

493      prospective study. Hum Reprod 2007;22:2455–62.

494   26.   Luke B, Brown MB, Wantman E, Stern JE, Baker VL, Widra E, et al. A prediction

495      model for live birth and multiple births within the first three cycles of assisted

496      reproductive technology. Fertil Steril 2014;102:744–52.

497   27.   McLernon DJ, Lee AJ, Maheshwari A, van Eekelen R, van Geloven N, Putter H, et al.

498      Predicting the chances of having a baby with or without treatment at different time

499      points in couples with unexplained subfertility. Hum Reprod 2019;34:1126–38.

500    28.    Metello JL, Tomás C, Ferreira P. Can we predict the IVF/ICSI live birth rate? J Bras

501            Reprod Assist 2019;23:402–7.

502    29.    Nelson SM, Fleming R, Gaudoin M, Choi B, Santo-Domingo K, Yao M.

503            Antimüllerian hormone levels and antral follicle count as prognostic indicators in a

504            personalized prediction model of live birth. Fertil Steril 2015;104:325–32.

505    30.    Nelson SM, Lawlor DA. Predicting live birth, preterm delivery, and low birth weight

506            in infants born from in vitro fertilisation: A prospective study of 144,018 treatment

507            cycles. PLoS Med 2011;8.

508    31.    Pettersson G, Nyboe Andersen A, Broberg P, Arce JC. Pre-stimulation parameters

509            predicting live birth after IVF in the long GnRH agonist protocol. Reprod Biomed

510            Online 2010;20:572–81.

511    32.    Porcu G, Lehert P, Colella C, Giorgetti C. Predicting live birth chances for women

512            with multiple consecutive failing IVF cycles: a simple and accurate prediction for

513            routine medical practice. Reprod Biol Endocrinol 2013;11:1.

514    33.    Ballester M, Oppenheimer A, d'Argent EM, Touboul C, Antoine J-M, Coutant C, et al.

515            Nomogram to predict pregnancy rate after ICSI-IVF cycle in patients with

516            endometriosis. Hum Reprod 2012;27:451–6.

517    34.    Qiu J, Li P, Dong M, Xin X, Tan J. Personalized prediction of live birth prior to the

518            first in vitro fertilization treatment: A machine learning method. J Transl Med 2019;17.

519    35.    Rongieres C, Colella C, Lehert P. To what extent does Anti-Mullerian Hormone

520            contribute to a better prediction of live birth after IVF? J Assist Reprod Genet

521            2015;32:37–43.

522    36.    Stolwijk AM, Straatman H, Zielhuis GA, Jansen CA, Braat DD, van Dop PA, et al.

523            External validation of prognostic models for ongoing pregnancy after in-vitro

524            fertilization. Hum Reprod 1998;13:3542–9.

23

525    37.    Templeton A, Morris JK, Parslow W. Factors that affect outcome of in-vitro

526            fertilisation treatment. Lancet 1996;348:1402–6.

527    38.    Wald M, Sparks AET, Sandlow J, Van-Voorhis B, Syrop CH, Niederberger CS.

528            Computational models for prediction of IVF/ICSI outcomes with surgically retrieved

529            spermatozoa. Reprod Biomed Online 2005;11:325–31.

530    39.    van Weert J-M, Repping S, van der Steeg JW, Steures P, van der Veen F, Mol BW. A

531            prediction model for ongoing pregnancy after in vitro fertilization in couples with male

532            subfertility. J Reprod Med 2008;53:250–6.

533    40.    Tarín JJ, Pascual E, García-Pérez MA, Gómez R, Hidalgo-Mora JJ, Cano A. A

534            predictive model for women's assisted fecundity before starting the first IVF/ICSI

535            treatment cycle. J Assist Reprod Genet 2020;

536    41.    Bancsi LFJMM, Huijs AM, Den Ouden CT, Broekmans FJM, Looman CWN,

537            Blankenstein MA, et al. Basal follicle-stimulating hormone levels are of limited value

538            in predicting ongoing pregnancy rates after in vitro fertilization. Fertil Steril

539            2000;73:552–7.

540    42.    Brodin T, Hadziosmanovic N, Berglund L, Olovsson M, Holte J. Comparing four

541            ovarian reserve markers--associations with ovarian response and live births after

542            assisted reproduction. Acta Obstet Gynecol Scand 2015;94:1056–63.

543    43.    Choi B, Bosch E, Lannon BM, Leveille MC, Wong WH, Leader A, et al. Personalized

544            prediction of first-cycle in vitro fertilization success. Fertil Steril 2013;99(7):1905–11.

545    44.    Dhillon RK, McLernon DJ, Smith PP, Fishel S, Dowell K, Deeks JJ, et al. Predicting

546            the chance of live birth for women undergoing IVF: A novel pretreatment counselling

547            tool. Hum Reprod 2016;31:84–92.

548    45.    Ferlitsch K, Sator MO, Gruber DM, Rücklinger E, Gruber CJ, Huber JC. Body mass

549            index, follicle-stimulating hormone and their predictive value in in vitro fertilization. J

550        Assist Reprod Genet 2004;21:431–6.

551  46.  Güvenir HA, Misirli G, Dilbaz S, Ozdegirmenci O, Demir B, Dilbaz B. Estimating the

552        chance of success in IVF treatment using a ranking algorithm. Med Biol Eng Comput

553        2015;53:911–20.

554  47.  Hamdine O, Eijkemans MJC, Lentjes EGW, Torrance HL, Macklon NS, Fauser

555        BCJM, et al. Antimüllerian hormone: prediction of cumulative live birth in

556        gonadotropin-releasing hormone antagonist treatment for in vitro fertilization. Fertil

557        Steril 2015;104:891-898.e2.

558  48.  Banerjee P, Choi B, Shahine LK, Jun SH, O'Leary K, Lathi RB, et al. Deep

559        phenotyping to predict live birth outcomes in in vitro fertilization. Proc Natl Acad Sci

560        U S A 2010;107:13570–5.

561  49.  Blank C, Wildeboer RR, DeCroo I, Tilleman K, Weyers B, de Sutter P, et al.

562        Prediction of implantation after blastocyst transfer in in vitro fertilization: a machine-

563        learning perspective. Fertil Steril 2019;

564  50.  Ho V, Pham T, Ho T, Vuong L. Predictive Model for Live Birth at 12 Months After

565        Starting In-Vitro Fertilization Treatment. MedPharmRes 2018;2:5–20.

566  51.  Hunault CC, Eijkemans MJC, Pieters MHEC, Te Velde ER, Habbema JDF, Fauser

567        BCJM, et al. A prediction model for selecting patients undergoing in vitro fertilization

568        for elective single embryo transfer. Fertil Steril 2002;77:725–32.

569  52.  Hunault CC, te Velde ER, Weima SM, Macklon NS, Eijkemans MJC, Klinkert ER, et

570        al. A case study of the applicability of a prediction model for the selection of patients

571        undergoing in vitro fertilization for single embryo transfer in another center. Fertil

572        Steril 2007;87:1314–21.

573  53.  Jones CA, Christensen AL, Salihu H, Carpenter W, Petrozzino J, Abrams E, et al.

574        Prediction of individual probabilities of livebirth and multiple birth events following in

575   vitro fertilization (IVF): a new outcomes counselling tool for IVF providers and

576   patients using HFEA metrics. J Exp Clin Assist Reprod 2011;8:3.

577   54.   Kaufmann SJ, Eastaugh JL, Snowden S, Smye SW, Sharma V. The application of

578   neural networks in predicting the outcome of in- vitro fertilization. Hum Reprod

579   1997;12:1454–7.

580   55.   Kim SK, Kim H, Oh S, Lee JR, Jee BC, Kim SH. Development of a novel nomogram

581   for predicting ongoing pregnancy after in vitro fertilization and embryo transfer.

582   Obstet Gynecol Sci 2018;61:669–74.

583   56.   Liao S, Xiong J, Tu H, Hu C, Pan W, Geng Y, et al. Prediction of in vitro fertilization

584   outcome at different antral follicle count thresholds combined with female age, female

585   cause of infertility, and ovarian response in a prospective cohort of 8269 women. Med

586   (United States) 2019;98.

587   57.   van Loendersloot LL, van Wely M, Repping S, Bossuyt PMM, van der Veen F.

588   Individualized decision-making in IVF: calculating the chances of pregnancy. Hum

589   Reprod 2013;28:2972–80.

590   58.   Meijerink AM, Cissen M, Mochtar MH, Fleischer K, Thoonen I, De Melker AA, et al.

591   Prediction model for live birth in ICSI using testicular extracted sperm. Adv Access

592   Publ July 2016;31:1942–51.

593   59.   Ottosen LDM, Kesmodel U, Hindkjær J, Ingerslev HJ. Pregnancy prediction models

594   and eSET criteria for IVF patients - Do we need more information? J Assist Reprod

595   Genet 2007;24:29–36.

596   60.   Cai QF, Wan F, Huang R, Zhang HW. Factors predicting the cumulative outcome of

597   IVF/ICSI treatment: a multivariable analysis of 2450 patients. Hum Reprod

598   2011;26:2532–40.

599   61.   Roberts SA, Fitzgerald CT, Brison DR. Modelling the impact of single embryo

600      transfer in a national health service IVF programme. Hum Reprod 2009;24(1):122–31.

601   62.  Roberts SA, Hirst WM, Brison DR, Vail A, towardSET collaboration. Embryo and

602      uterine influences on IVF outcomes: an analysis of a UK multi-centre cohort. Hum

603      Reprod 2010;25:2792–802.

604   63.  Roberts SA, Hann M, Brison DR. Factors affecting embryo viability and uterine

605      receptivity: insights from an analysis of the UK registry data. Reprod Biomed Online

606      2016;32:197–206.

607   64.  Sunkara SK, Rittenberg V, Raine-Fenning N, Bhattacharya S, Zamora J,

608      Coomarasamy A. Association between the number of eggs and live birth in IVF

609      treatment: an analysis of 400 135 treatment cycles. Hum Reprod 2011;26:1768–74.

610   65.  Uyar A, Bener A, Ciray HN. Predictive Modeling of Implantation Outcome in an in

611      Vitro Fertilization Setting. Med Decis Mak 2015;35:714–25.

612   66.  Verberg MFG, Eijkemans MJC, Macklon NS, Heijnen EMEW, Fauser BCJM,

613      Broekmans FJ. Predictors of ongoing pregnancy after single-embryo transfer following

614      mild ovarian stimulation for IVF. Fertil Steril 2008;

615   67.  Vaegter KK, Lakic TG, Olovsson M, Berglund L, Brodin T, Holte J. Which factors are

616      most predictive for live birth after in vitro fertilization and intracytoplasmic sperm

617      injection (IVF/ICSI) treatments? Analysis of 100 prospectively recorded variables in

618      8,400 IVF/ICSI single-embryo transfers. Fertil Steril 2017;107:641-648.e2.

619   68.  Vaegter KK, Berglund L, Tilly J, Hadziosmanovic N, Brodin T, Holte J. Construction

620      and validation of a prediction model to minimize twin rates at preserved high live birth

621      rates after IVF. Reprod Biomed Online 2019;38:22–9.

622   69.  Vogiatzi P, Pouliakis A, Siristatidis C. An artificial neural network for the prediction

623      of assisted reproduction outcome. J Assist Reprod Genet 2019;36:1441–8.

624   70.  Wu F, Liu F, Guan Y, Du J, Tan J, Lv H, et al. A nomogram predicting clinical

625        pregnancy in the first fresh embryo transfer for women undergoing in vitro fertilization

626        and intracytoplasmic sperm injection (IVF/ICSI) treatments. J Biomed Res

627        2019;33:422.

628    71.    Carrera-Rotllan J, Estrada-García L, Sarquella-Ventura J. Prediction of pregnancy in

629        IVF cycles on the fourth day of ovarian stimulation. J Assist Reprod Genet

630        2007;24:387–94.

631    72.    Tarín JJ, Pascual E, Gómez R, García-Pérez MA, Cano A. Predictors of live birth in

632        women with a history of biochemical pregnancies after assisted reproduction

633        treatment. Eur J Obstet Gynecol Reprod Biol 2020;

634    73.    Corani G, Magli C, Giusti A, Gianaroli L, Gambardella LM. A Bayesian network

635        model for predicting pregnancy after in vitro fertilization. Comput Biol Med

636        2013;43:1783–92.

637    74.    Dessolle L, Fréour T, Ravel C, Jean M, Colombel A, Daraï E, et al. Predictive factors

638        of healthy term birth after single blastocyst transfer. Hum Reprod 2011;

639    75.    Gianaroli L, Magli MC, Gambardella L, Giusti A, Grugnetti C, Corani G. Objective

640        way to support embryo transfer: a probabilistic decision. Hum Reprod 2013;28:1210–

641        20.

642    76.    Goldman RH, Kaser DJ, Missmer SA, Srouji SS, Farland L V, Racowsky C. Building

643        a model to increase live birth rate through patient-specific optimization of embryo

644        transfer day. J Assist Reprod Genet 2016;33:1525–32.

645    77.    Grin L, Mizrachi Y, Cohen O, Lazer T, Liberty G, Meltcer S, et al. Does progesterone

646        to oocyte index have a predictive value for IVF outcome? A retrospective cohort and

647        review of the literature. Gynecol Endocrinol 2018;34:638–43.

648    78.    Hirst WM, Vail A, Brison DR, Roberts SA. Prognostic factors influencing fresh and

649        frozen IVF outcomes: an analysis of the UK national database. Reprod Biomed Online

650    2011;22:437–48.

651    79.    McLernon DJ, Steyerberg EW, Te Velde ER, Lee AJ, Bhattacharya S. Predicting the

652            chances of a live birth after one or more complete cycles of in vitro fertilisation:

653            Population based study of linked cycle data from 113 873 women. BMJ

654            2016;355:i5735.

655    80.    Leijdekkers JA, Eijkemans MJC, van Tilborg TC, Oudshoorn SC, McLernon DJ,

656            Bhattacharya S, et al. Predicting the cumulative chance of live birth over multiple

657            complete cycles of in vitro fertilization: an external validation study. Hum Reprod

658            2018;33:1684–95.

659    81.    Stolwijk AM, Zielhuis GA, Hamilton CJCM, Straatman H, Hollanders JMG, Goverde

660            HJM, et al. Prognostic models for the probability of achieving an ongoing pregnancy

661            after in-vitro fertilization and the importance of testing their predictive value. Hum

662            Reprod 1996;11:2298–303.

663    82.    Khader A, Lloyd SM, McConnachie A, Fleming R, Grisendi V, La Marca A, et al.

664            External validation of anti-Müllerian hormone based prediction of live birth in assisted

665            conception. J Ovarian Res 2013;6:3.

666    83.    Qiu J, Li P, Dong M, Xin X, Tan J. Live birth prediction before the first IVF treatment.

667    84.    University of Aberdeen. Outcome Prediction in Subfertility.

668    85.    Society for Assisted Reproductive Technology. What are my chances with ART?

669            [Internet]. 2020;Available from: https://www.sartcorsonline.com/Predictor/Patient

670    86.    Choi B, Bosch E, Lannon BM, Leveille M-C, Wong WH, Leader A, et al. Personalized

671            prediction of first-cycle in vitro fertilization success. Fertil Steril 2013;99:1905–11.

672    87.    Sperrin M, Martin GP, Pate A, Van Staa T, Peek N, Buchan I. Using marginal

673            structural models to adjust for treatment drop-in when developing clinical prediction

674            models. Stat Med 2018;

675  88.  Hu YH, Wu F, Lo CL, Tai CT. Predicting warfarin dosage from clinical data: A

676       supervised learning approach. Artif Intell Med 2012;

677  89.  Frank I, Blute ML, Cheville JC, Lohse CM, Weaver AL, Zincke H. An outcome

678       prediction model for patients with clear cell renal cell carcinoma treated with radical

679       nephrectomy based on tumor stage, size, grade and necrosis: The SSIGN score. J Urol

680       2002;

681  90.  Kaliarnta S, Nihlén-Fahlquist J, Roeser S. Emotions and ethical considerations of

682       women undergoing IVF-treatments. HEC Forum 2011;

683  91.  Mol BW, Verhagen TEM, Hendriks DJ, Collins JA, Coomarasamy A, Opmeer BC, et

684       al. Value of ovarian reserve testing before IVF: A clinical decision analysis. Hum

685       Reprod 2006;

686  92.  Nachtigall RD, Dougall K Mac, Lee M, Harrington J, Becker G. What do patients

687       want? Expectations and perceptions of IVF clinic information and support regarding

688       frozen embryo disposition. Fertil Steril 2010;

689  93.  Vollmer S, Mateen BA, Bohner G, Király FJ, Ghani R, Jonsson P, et al. Machine

690       learning and artificial intelligence research for patient benefit: 20 critical questions on

691       transparency, replicability, ethics, and effectiveness. BMJ 2020;

692  94.  Hassan MR, Al-Insaif S, Hossain MI, Kamruzzaman J. A machine learning approach

693       for prediction of pregnancy outcome following IVF treatment. Neural Comput Appl

694       2020;

695  95.  Chen JH, Asch SM. Machine learning and prediction in medicine-beyond the peak of

696       inflated expectations. N. Engl. J. Med. 2017;376:2507–9.

697

698

699    **Figure legends:**

700

701    **Figure (1):** Study selection and inclusion process on prediction models for reproductive

702    outcomes following assisted reproductive technology treatments.

703

704    **Figure (2):** TRIPOD classification of included studies reporting on prediction models for

705    reproductive outcomes following assisted reproductive technology treatments

706

707    **Figure (3):** Predictors used in the development of prediction models for reproductive

708    outcomes following assisted reproductive technology treatments.

709    3a: predictors in pre-ART treatment models

710    3b: predictors for intra-ART treatment models
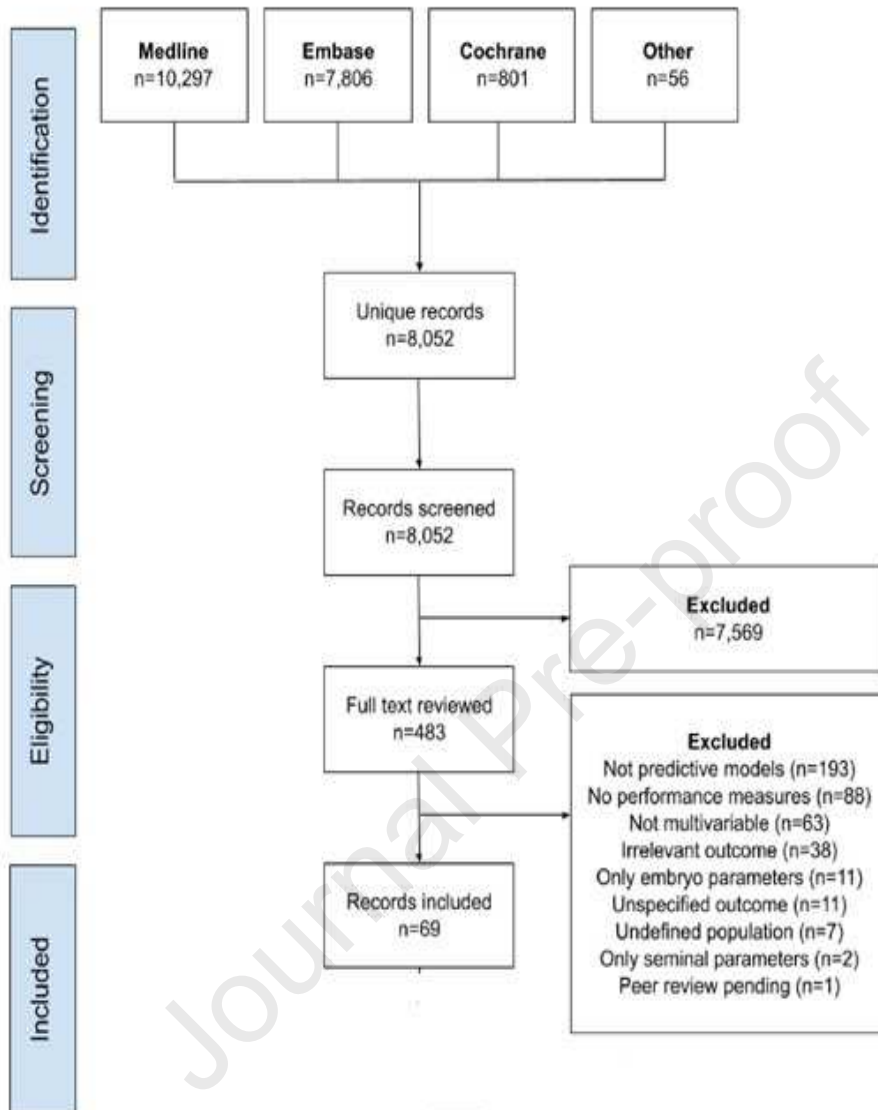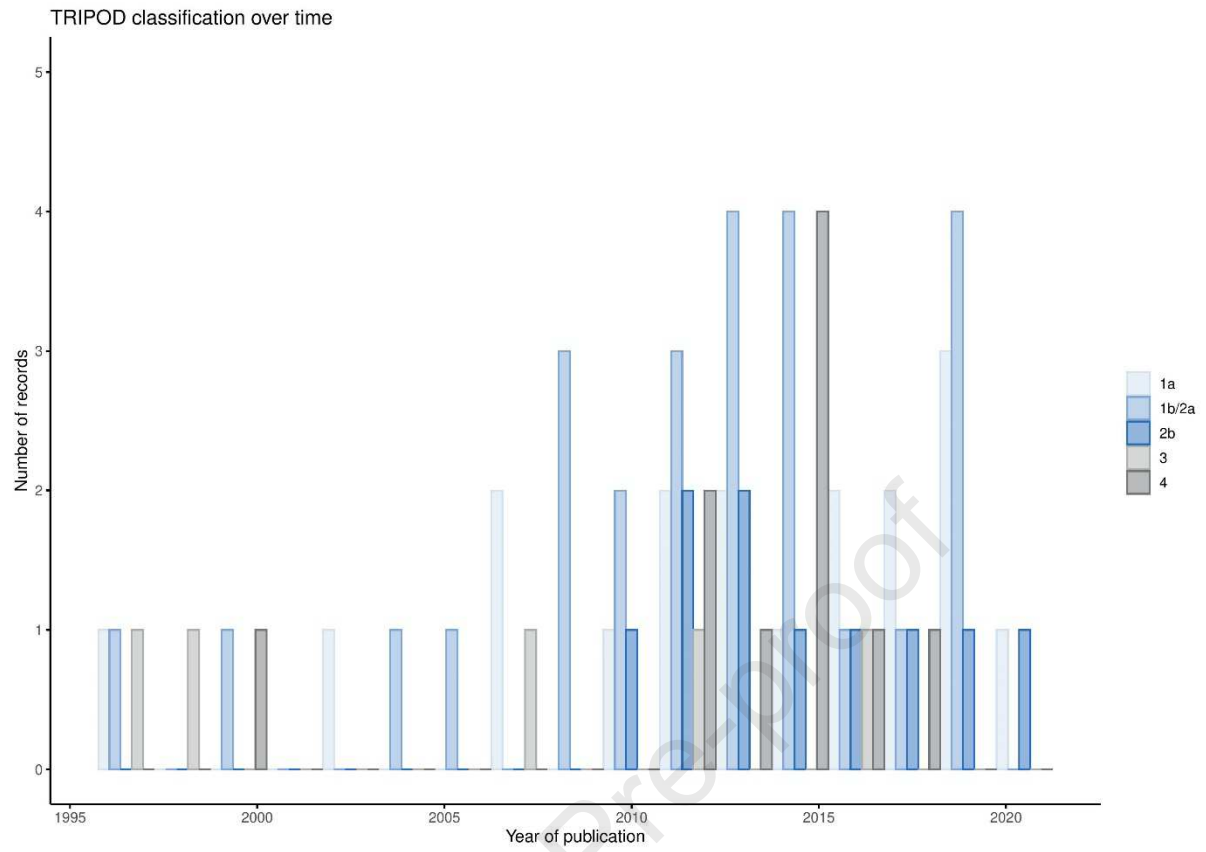
711

712    **Figure (4):** Risk of bias assessment in included studies reporting on prediction models for

713    reproductive outcomes following assisted reproductive technology treatments
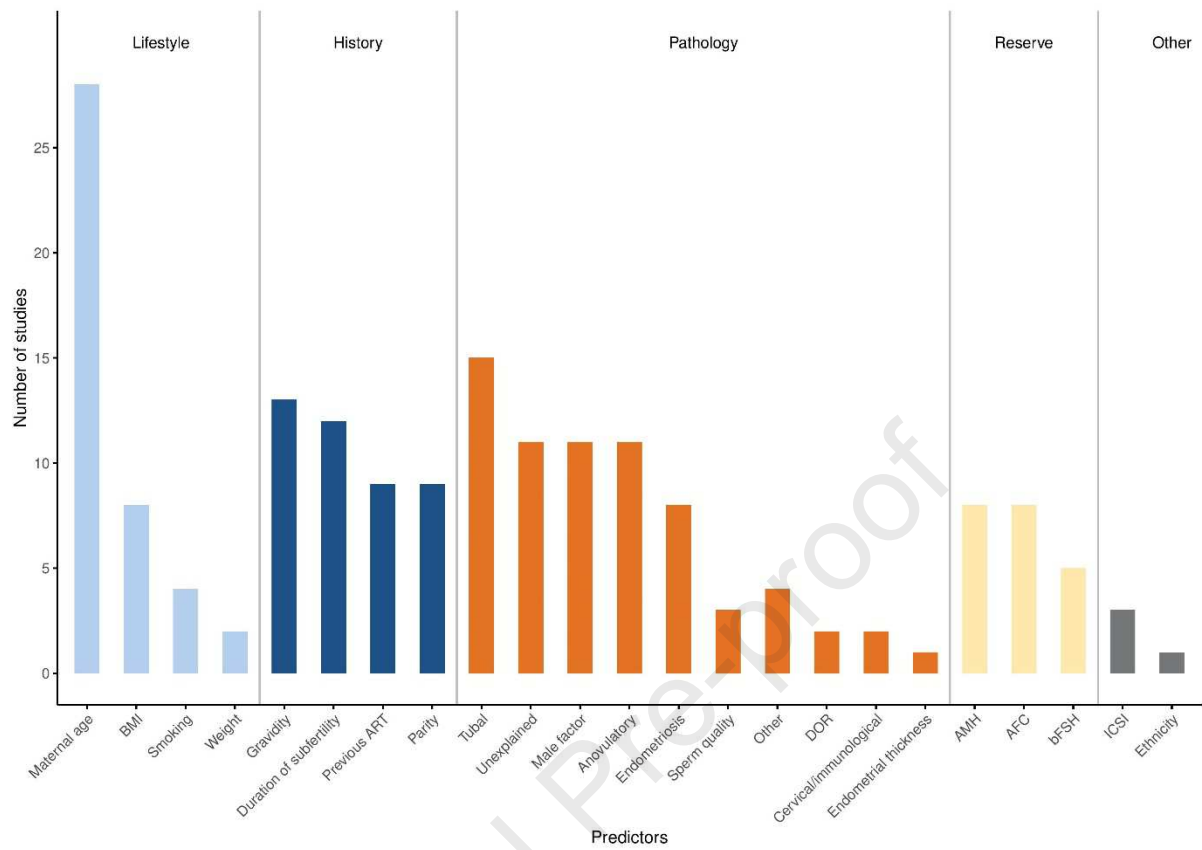
714
715

TRIPOD classification over time

**3a**



**3b**