

1 **Title**

2 Integrating expert opinions with clinical trial data to analyse low-powered subgroup analyses: a
3 Bayesian analysis of the VerDiCT trial

4 **Authors:**

5 Russell Thirard ¹, Raimondo Ascione ², Jane M Blazeby ^{3,4} and Chris A Rogers ^{1,3}

6 **Affiliations:**

7 ¹ Bristol Trials Centre (BTC), University of Bristol, Bristol, UK

8 ² Bristol Medical School, Bristol Heart Institute, University of Bristol, Bristol, UK

9 ³ National Institute for Health Research Bristol Biomedical Research Centre, Bristol, UK

10 ⁴ Division of Surgery, University Hospitals Bristol NHS Foundation Trust, Bristol, UK

11

12 **Corresponding author:**

13 Russell Thirard

14 Bristol Trials Centre, Bristol Medical School

15 University of Bristol

16 Zone A, level 7, Bristol Royal Infirmary

17 Bristol BS2 8HW, UK

18 Email: russell.thirard@bristol.ac.uk

19

20

21

22 Abstract

23 Background

24 Typically, subgroup analyses in clinical trials are conducted by comparing the intervention
25 effect in each subgroup by means of an interaction test. However, trials are rarely, if ever,
26 adequately powered for interaction tests, so clinically important interactions may go
27 undetected. We discuss the application of Bayesian methods by using expert opinions
28 alongside the trial data. We applied this methodology to the VerDiCT trial investigating the
29 effect of preoperative volume replacement therapy (VRT) versus no VRT (usual care) in
30 diabetic patients undergoing cardiac surgery. Two subgroup effects were of clinical interest,
31 a) preoperative renal failure and b) preoperative type of antidiabetic medication.

32 Methods

33 Clinical experts were identified within the VerDiCT trial centre in the UK. A questionnaire
34 was designed to elicit opinions on the impact of VRT on the primary outcome of time from
35 surgery until medically fit for hospital discharge, in the different subgroups. Prior beliefs of
36 the subgroup effect of VRT were elicited face-to-face using two unconditional and one
37 conditional questions per subgroup analysis. The robustness of results to the 'community of
38 priors' was assessed. The community of priors was built using the expert priors for the mean
39 average treatment effect, the interaction effect or both in a Bayesian Cox proportional
40 hazards model implemented in the STAN software in R.

41 Results

42 Expert opinions were obtained from 7 clinicians (6 cardiac surgeons and 1 cardiac
43 anaesthetist). Participating experts believed VRT could reduce the length of recovery

44 compared to usual care and the greatest benefit was expected in the subgroups with the
45 more severe comorbidity. The Bayesian posterior estimates were more precise compared to
46 the frequentist maximum likelihood estimate and were shifted toward the overall mean
47 treatment effect.

48 Conclusions

49 In the VeRDICT trial, the Bayesian analysis did not provide evidence of a difference in
50 treatment effect across subgroups. However, this approach increased the precision of the
51 estimated subgroup effects and produced more stable treatment effect point estimates
52 than the frequentist approach. Trial methodologists are encouraged to prospectively
53 consider Bayesian subgroup analyses when low-powered interaction tests are planned.

54 Trial registration

55 ISRCTN, ISRCTN02159606. Registered 29th October 2008.

56 **Keywords:** Bayesian analysis, Subgroup analyses, Survival, Cardiac surgery, Volume
57 replacement therapy, Elicitation

58

59 Introduction

60 In clinical trials, evaluation of intervention effects in subgroups of patients may be pre-
61 specified in the study protocol, or exploratory analyses carried out post-hoc. The aim of
62 such analyses is to assess if the intervention effect is consistent across patients or if specific
63 subgroups of patients experience larger benefit or harm. Typically, subgroup analyses in
64 clinical trials are conducted by comparing the intervention effect in each subgroup by means
65 of an interaction test in the frequentist framework. In trials reporting subgroup analyses,

66 between 27 and 34.5% planned a statistical test for interaction [1-3]. When the test does
67 not provide evidence of an interaction, it is recommended to report the overall results.
68 However, these tests are low-powered and clinically important interactions in subgroup of
69 patients may go undetected. In light of this, some suggest raising the Type I error rate when
70 testing interactions, thereby increasing power but this is offset by the increased risk of
71 “false positives” and chance findings [4].

72 A Bayesian analysis differs from the frequentist analysis in that the uncertainty about
73 unknown parameters such as an interaction parameter can be expressed in a ‘prior’
74 distribution. Instead of assuming that there is no prior knowledge about the interaction
75 effect, or assuming that there is no interaction by pooling all patients together irrespective
76 of their subgroup, the proposed Bayesian approach allows the user to flexibly apply an
77 expert informative prior for the interaction parameter using external knowledge to the trial.
78 Data from previous studies and/or expert opinions can be used in a prior to characterise the
79 possibly differential effects of an intervention for subgroups of patients having different
80 comorbidities. By Bayes’ theorem, inferences about the subgroup treatment effects are
81 drawn from the ‘posterior’ distributions. The posterior distributions are proportional to the
82 product of the expert prior distributions and the trial data also referred as the ‘likelihood’.
83 The interest of this method lies in situations where increasing the power by collecting more
84 patient-level data is not feasible, too expensive or too time-consuming.

85 **Motivating case study**

86 The VeRDICT study comprised two randomised controlled trials conducted in parallel. The
87 study investigated the effect of volume replacement therapy (VRT) on postoperative length
88 of stay before being fit-for-discharge. The population of interest were diabetic patients

89 undergoing coronary artery bypass grafting surgery. The study randomised 169 patients
90 (122 in the UK and 47 in India). Two subgroup effects were of clinical interest, a)
91 preoperative renal function and b) preoperative blood glucose management. In relation to
92 subgroup a), VRT effect was compared in patients with evidence of either microalbuminuria
93 or diabetic nephropathy to patients without (low risk of renal failure subgroup A1 versus
94 high risk subgroup A2). With respect to subgroup b), the analysis compared patients
95 managing their diabetes only with oral medication (less severe comorbidity subgroup B1) to
96 patients taking only insulin or combined with oral medication (more severe comorbidity
97 subgroup B2). One hundred and seventy participants were required in order to detect a
98 significant intervention effect in the combined trial population using frequentist methods.
99 Recruiting as many in each subgroup was not feasible as the intervention is very specialised
100 and would have extended the recruitment period by more than four years of recruitment.
101 Hence, the aim of this study was to apply a Bayesian approach to analyse the subgroup
102 effects in the VeRDICT trial and examine if this method was valuable.

103 **Methods**

104 **Prior elicitation and expert opinion derivation**

105 An extensive literature is available on elicitation methods and associated heuristics that may
106 bias an expert's ability to assess probabilities [5-8]. Examples of elicitation questionnaires in
107 medical trials have been published [9-12]. Our questionnaire was based on the work of
108 Spiegelhalter et al. [13]. Experts were identified as health professionals well informed about
109 VRT and effects in diabetic patients undergoing CABG surgery. Cardiac surgeons and cardiac
110 anaesthetists were identified within the UK trial site. The questionnaire was delivered face-
111 to-face and as part of a routine research meeting before the results of the VeRDICT trial

112 were available. As hazards are complex to elicit, experts were asked about the number of
113 patients they would expect to be fit-for-discharge within six days post-operation. Six days
114 post-operation was determined as the median length of stay for diabetic patients
115 undergoing CABG surgery between 2012 and 2015 from the Patient Analysis and Tracking
116 System (PATS) database – a registry of every adult cardiac surgical procedure undertaken at
117 the Bristol Royal Infirmary cardiac surgical unit. Being fit-for-discharge within six days was
118 defined as the outcome in the questionnaire as a ‘normal’ recovery without any major
119 complications.

120 The experts were asked to express their opinions in the quantile format as suggested by
121 Cooke [14, 15]. Experts provided percentiles of their subjective distribution using the
122 median, 2.5 and 97.5 percentiles (i.e. the most likely value, the lowest and highest plausible
123 values respectively). Experts were encouraged to imagine as if they would have to provide
124 bounds of a 95% range of plausible values (see supplementary material).

125 The questionnaire answers were used as prior beliefs for the hazard ratios. Hence a
126 transformation was required to convert the opinions to the log hazard ratio scale:

129
$$\log(HR) = \log \left[\frac{\log (1 - p_1)}{\log (1 - p_2)} \right]$$

127 where p_1 and p_2 are the probabilities of having a normal recovery in the VRT group and
128 usual care group respectively.

130 A positive $\log(HR)$ (i.e. $HR > 1$) suggests that the ‘hazard’ of having a normal recovery is
131 higher and therefore, favours the VRT group compared to the usual care group. The
132 $\log(HR)$ was assumed to follow an approximate normal likelihood distribution. Each
133 expert’s elicited opinions were fitted with normal distributions using least squares on the
134 cumulative distribution function [16]. The experts’ individual distributions were then

135 aggregated to obtain a combined experts prior distribution using the linear opinion pooling
136 method [17].

137 The intervention effects on the log hazard ratio scale are denoted θ_1 and θ_2 for the less and
138 more severe subgroups respectively such that $\theta_1 \sim N(p_1, s_1^2)$ and $\theta_2 \sim N(p_2, s_2^2)$. Experts
139 were firstly asked to provide their opinions about the number of patients treated with VRT
140 having a normal recovery in the less and more severe subgroups separately in each
141 unconditional question. This was then followed by a conditional question which asked the
142 experts how their opinions changed about the effect in the more severe subgroups given
143 two scenarios in the less severe subgroups. The two scenarios were: a) if we knew the true
144 effect in the less severe subgroup is null (i.e. $\theta_1 = 0$); and b) if we knew the true effect in
145 the less severe subgroup is beneficial by a value d (i.e. $\theta_1 = d, d > 0$). The following
146 distributions were derived: $\theta_2 | (\theta_1 = 0) \sim N(p_{20}, s_{20}^2)$ and $\theta_2 | (\theta_1 = d) \sim N(p_{2d}, s_{2d}^2)$. For
147 the risk of renal failure subgroup analysis, the beneficial change was $d = \log \left[\frac{\log(1-70/100)}{\log(1-60/100)} \right]$
148 meaning we supposed that in truth VRT increased the number of patients with low risk of
149 renal failure having a normal recovery to 70 compared to 60 when treated with usual care.
150 For the blood glucose management subgroup analysis, the change was

151 $d = \log \left[\frac{\log(1-65/100)}{\log(1-50/100)} \right]$. The joint prior distribution of the treatment effect in each subgroup

152 $\begin{pmatrix} \theta_1 \\ \theta_2 \end{pmatrix}$ was modelled as a bivariate Normal with mean $m = (p_1, p_2)$ and variance-covariance

153 matrix $V = \begin{pmatrix} V_{11} & V_{12} \\ V_{12} & V_{22} \end{pmatrix}$. The parameter derivations followed the methods outlined by White

154 et al. [18]:

- 155 1. The unconditional variance V_{11} for the effect of VRT in the less severe subgroup was
156 derived from the variance of the elicited distribution: $V_{11} = s_1^2$.

- 157 2. The covariance element V_{12} was derived as $V_{12} = b_{12}V_{11}$. The regression coefficient
 158 (b_{12}) of the treatment effect in the more severe subgroup on the treatment effect in
 159 the less severe subgroup is defined as $b_{12} = \frac{p_{2d}-p_{20}}{d-0}$.
- 160 3. The variance of the treatment effect in the more severe subgroup V_{22} is derived such
 161 that $V_{22} = var(\theta_2|\theta_1) + b_{12}^2 V_{11}$. The variance of the treatment effect in the more
 162 severe subgroup conditional on the treatment effect in the less severe subgroup was
 163 derived as the average of the two conditional question variances:
 164 $var(\theta_2|\theta_1) = \frac{s_{20}^2+s_{2d}^2}{2}$. A sensitivity analysis assuming the variance $V_{22} = s_2^2$ was
 165 performed to compare posterior estimates to alternative variance derivations when
 166 the resulting variance-covariance was positive definite.
- 167 4. A further sensitivity analysis was performed by drawing the mean of θ_2 from the
 168 regression $E(\theta_2|\theta_1) = p_{20} + b_{12}\theta_1$, giving $E(\theta_2) = p_{20} + b_{12}p_1$.

169 Community of prior specifications

170 This paper assessed results to a wide variety of expert opinions and ways of using the expert
 171 opinions. A clinical prior, a sceptical prior, an interaction prior, and other specifications of
 172 the expert opinions composed our 'community of priors' and are described below:

- 173 • The clinical prior was directly derived from the experts' opinions: $\theta \sim N_2(m, V)$
- 174 • The sceptical prior used the experts' opinions for the variance component, but the
 175 mean was centred around the null effect: $\theta \sim N_2(0, V)$
- 176 • The interaction prior only used an informative expert prior for the treatment-by-
 177 subgroup interactions. This was enabled by re-parameterising θ such that the prior ψ
 178 consisted of an uninformative prior for the mean average VRT effect and an informative

179 prior for the interaction term comparing the effect in the more severe subgroup to that in
180 the less severe subgroup:

$$181 \quad \psi = C \theta = \begin{pmatrix} 1/2 & 1/2 \\ -1 & 1 \end{pmatrix} \times \begin{pmatrix} \theta_1 \\ \theta_2 \end{pmatrix} = \begin{pmatrix} \frac{1}{2}(\theta_1 + \theta_2) \\ \theta_2 - \theta_1 \end{pmatrix}$$

182 The prior ψ variance is $W = CV C^T$. The matrix W was modified such that:

183 $W^* = \begin{pmatrix} L & 0 \\ 0 & W_{22} \end{pmatrix}$ with L an arbitrarily large number. L in W^* expresses the lack of
184 information about the average treatment effect whilst W_{22} expresses the prior information
185 from our experts about the interaction effect $\theta_2 - \theta_1$. The off-diagonal elements express
186 the absence of any covariance beliefs between the average intervention effect and the
187 interaction effect. A matrix V^* was back transformed such that $V^* = C^{-1}W^*C^{-T}$. The
188 interaction prior followed $\theta \sim N_2(m, V^*)$. The Bayesian analysis using the interaction prior
189 specification was defined as the primary analysis of the VerDiCT trial.

190 • A further specification of the interaction prior $\theta \sim N_2(0, V^*)$ was assessed to avoid
191 reporting qualitative interactions (e.g. $\theta_1 < 0$ and $\theta_2 > 0$) in specific cases where the
192 treatment effects were null in both subgroups but the experts believed VRT was much
193 better in the more severe subgroup than in the less severe subgroup (i.e. $\theta_2 - \theta_1 > 0$).

194 • A vague prior that would closely reproduce the results of a frequentist maximum
195 likelihood estimate (MLE): $\theta \sim N_2(0, LI)$ where L is an arbitrarily large number and I is the
196 identity matrix. This prior specification was added following comments from clinical experts
197 asking about the performance of the Bayesian approach when there is no knowledge about
198 the subgroup effects. Table 1 presents a summary of the community of priors and their
199 associated mean and variance parameters.

200 *Table 1 Summary of prior specifications*

Prior specifications	Mean parameter	Variance parameters	
		Mean effect	Interaction effect
Clinical	m	V ^a	V ^a
Sceptical	0	V ^a	V ^a
Interaction	m	V ^{*b}	V ^{*a}
Interaction-variance	0	V ^{*b}	V ^{*a}
Vague	0	LI ^b	LI ^b

201 ^a Informative using expert opinions ^b Uninformative

202 **Statistical analysis**

203 The VeRDICT primary outcome was time-to-event and was analysed using an adjusted Cox
 204 proportional hazards model in the trial report [19]. A Bayesian counterpart of the Cox
 205 proportional hazard model was performed adjusting for the same covariates. Analyses were
 206 performed using STAN software from R [20, 21] and STATA version 15.1 (StataCorp LP,
 207 College Station, TX, USA). White et al. proposed a normal-normal conjugate Bayesian analysis
 208 as the number of events was large and therefore the log-likelihood was reasonably
 209 approximated by a normal distribution [18]. However, we did not take this approach but
 210 jointly estimated the log-likelihood and the prior in a fully Bayesian Cox model to extend the
 211 proposed method to scenarios of low number of events. The STAN code for this analysis is
 212 provided in the supplementary materials [see Additional file 1]. As the experts were all from
 213 the UK site, the analysis was restricted to the UK trial. The results of the Bayesian analysis
 214 were compared with the frequentist approach to assess relative merits of each statistical
 215 framework.

216 **Results**

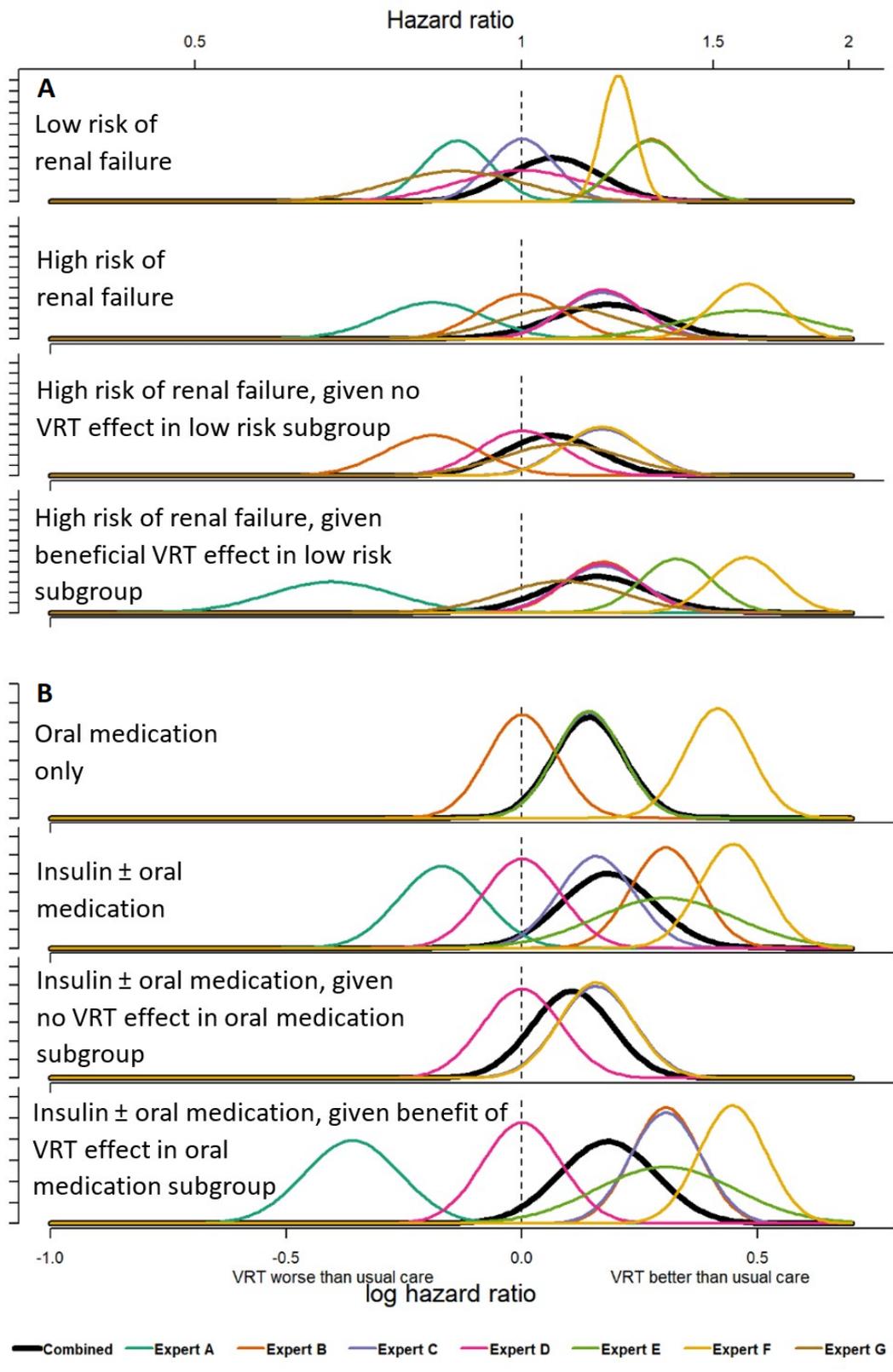
217 ***VeRDICT trial and expert results***

218 *Table 2 Baseline characteristics of the UK VeRDICT trial participants*

Treatment allocation	Randomised to usual care (n=61) n (%)	Randomised to VRT (n=60) n (%)
<i>Minimisation criteria</i>		
Age >70 years	17 (28)	17 (28)
Female gender	10 (16)	10 (17)
Preoperative creatinine >160 µmol/L	3 (5)	4 (7)
Ejection fraction <50%	49 (80)	49 (82)
Cardiac angiogram in the 5 days prior to surgery	5 (8)	7 (12)
<i>Subgroup variables</i>		
Low risk of renal failure	37 (62) ¹	38 (63)
Oral diabetic medication only	30 (49)	35 (58)

219 ¹ 1 patient with missing data

220 Table 2 presents baseline characteristics of the 121 participants who underwent
 221 randomisation in the UK trial of the VeRDICT study. The baseline characteristics were similar
 222 across treatment groups. The pre-surgery blood glucose management presents a slight
 223 imbalance with a higher proportion of patients managing their blood glucose with oral
 224 medication in this those randomised to VRT compared to patients randomised to usual care.



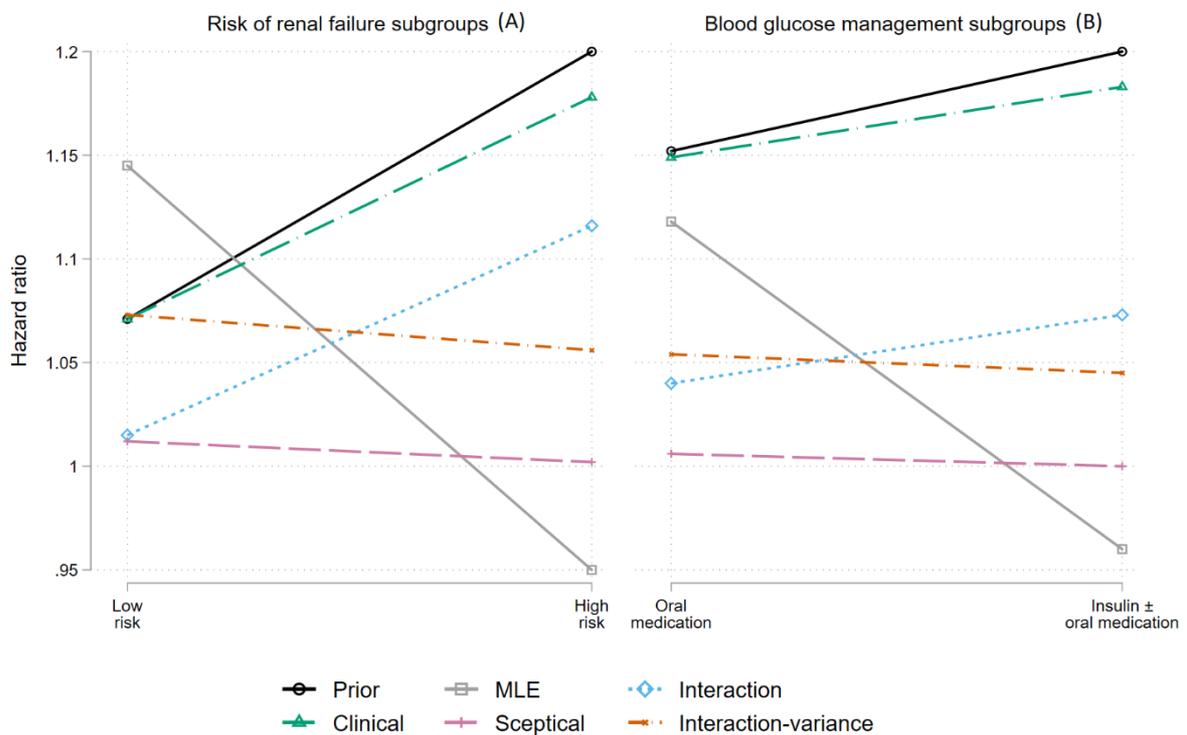
225

226 *Figure 1 Elicited prior distributions of the log hazard ratio for the risk of renal failure*
 227 *subgroups (A) and for the type of antidiabetic medication subgroups (B)*

228 Apart from the opinion of expert A, the unanimous belief was that VRT was better or at
 229 least not worse than usual care for all subgroups of patients. Additionally, the experts
 230 believed that VRT would provide the greatest benefit for higher risk patients with more
 231 severe comorbidities. For example, patients with high risk of renal failure or patients treated
 232 with insulin ± oral medication were expected to have 20% higher ‘hazard’ of a normal
 233 recovery with VRT compared to usual care (HR=1.20, 95% CI 0.95-1.51 and HR=1.20, 95% CI
 234 0.99-1.46 respectively). Nonetheless, the combined experts’ prior (presented in black) still
 235 includes the null effect suggesting there is some uncertainty about VRT being better than
 236 usual care.

237

238 **Bayesian analysis results**



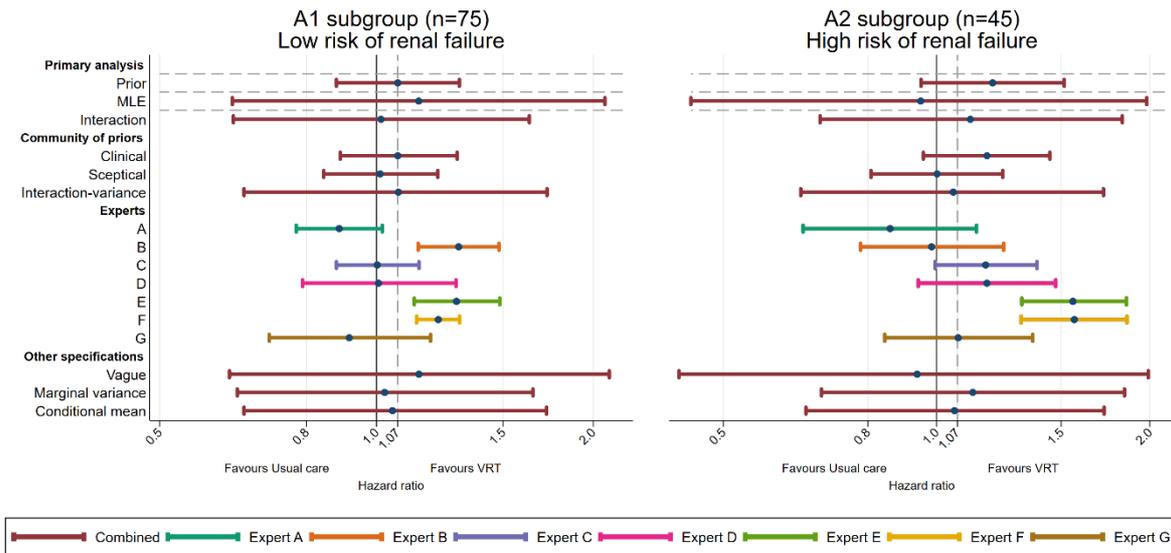
239
 240 *Figure 2 Expert prior median and posterior median estimates using the pooled expert prior*
 241 *with different specifications.*

242 The results of the Bayesian analyses combining the trial data and the expert opinions were
 243 very similar for both subgroup analyses A and B: the trial data treatment estimates (MLEs)

244 and the expert opinions (prior medians) were of the comparable magnitude and direction in
245 the risk of renal failure analysis (left panel) compared to the blood glucose management
246 subgroup analysis (right panel) in Figure 2. Hence for the following results, we have focused
247 on the risk of renal failure subgroup analysis.

248 While the experts believed that VRT would be of greater benefit in the higher risk subgroup
249 (shown in black), this was not borne out of the trial data (grey); the maximum likelihood
250 hazard ratio was lower in the high-risk group than in the low-risk group. The posterior
251 results using an expert prior for the interaction effect but a vague prior for the mean effect
252 (i.e. the interaction prior shown in blue) produces a slope that is similar to that of the prior
253 but shifted towards the trial's mean treatment effect (HR=1.07). The clinical posterior
254 estimates are closely aligned to the prior, in location and in slope, as it uses the expert
255 beliefs for both the overall mean and the interaction effects. This result suggests, the clinical
256 prior dominates the trial data when using informative priors for both the interaction effect
257 and the mean effect. The slopes of the posterior estimates using the sceptical and the
258 interaction-variance priors lie between the MLE and the prior, with their location centred
259 around the null effect and the mean treatment effect respectively.

260



261
 262 *Figure 3 Prior median and 95% Interval, MLE estimate and 95% confidence intervals, and*
 263 *posterior medians and 95% credible intervals for the treatment effect in the risk of renal*
 264 *failure subgroups*

265 Figure 3 presents the prior median, MLE and posteriors estimates for the risk of renal failure
 266 subgroup analysis and their variability (see supplementary material for blood glucose
 267 management analysis results). In this trial, the 95% confidence intervals for MLE are wide
 268 (e.g. 95% CI 0.63-2.08 and 95% CI 0.45-1.98 for A1 and A2 subgroups respectively) due to
 269 small subsamples (75 patients in subgroup A1, 45 in A2, 65 in B1 and 56 in B2). The
 270 variability of the combined experts' prior distribution was lower than the variability in the
 271 trial treatment effect MLEs.

272 The Bayesian posterior results using expert beliefs have increased precision compared to
 273 the MLE. In the presence of a qualitative interaction (i.e. the treatment effects in the
 274 subgroups are in opposite directions) which was observed in the trial data MLEs, the
 275 posterior estimates using the interaction prior are pulled towards the mean treatment
 276 effect. The precision of the estimates is increased but the results do not provide enough
 277 evidence of a benefit of VRT (52% and 67% probability of VRT being better than usual care
 278 for less and more severe subgroup respectively). Using each expert prior individually gave

279 posterior estimates which varied across experts. The posterior estimates using expert A's
280 prior, who was the most sceptical about VRT, favours usual care, with probabilities of only
281 5% and 15% of benefit with VRT in the less and more severe subgroups respectively. In
282 contrast, experts E and F were very enthusiastic about the effect of VRT and the posterior
283 estimate reports a 100% probability of VRT being better for both subgroups.

284 The posterior estimates using the vague prior confirms that when using uninformative priors
285 for all the parameters, the results are consistent with the MLEs that only use trial data for
286 inferences. Although the posterior results were robust to the different variance derivation,
287 the posterior results using the conditional mean derivation returned point estimates that
288 were different compared to the unconditional mean derivation (HR=1.05 vs HR=1.01 in the
289 low risk subgroup and HR=1.06 vs HR=1.12 in the high risk subgroup). We argue the prior
290 mean of θ_2 derived from the experts' answers to the second unconditional question is a
291 direct derivation requiring less manipulation than the conditional mean derivation and
292 therefore the unconditional method should be preferred. The posteriors results estimates
293 are provided in the supplementary files [see Additional files 2 and 3].

294 Discussion

295 A Bayesian approach to subgroup analyses has been successfully applied to the UK VeRDICT
296 trial. Expert opinions were elicited without knowledge of the trial results and these opinions
297 were combined with the trial data under a range of prior specifications. The questionnaire
298 eliciting opinions was challenging to design as it needed to not be overburdensome but at
299 the same time capture the information required (i.e. personal opinions on the effect of VRT
300 in each subgroup and the interaction effect between the intervention and the subgroups).
301 Face-to-face elicitation had the advantage that training and feedback could be provided to

302 the experts iteratively. Other elicitation methods (e.g. postal questionnaires) were
303 considered but we chose a face-to-face approach to allow us the opportunity to answer
304 questions and clarify what was being asked where required. As previously reported [22-24],
305 experts experienced challenges when characterising a range of plausible values and the
306 elicitation range may influence how they express their opinion.

307 The Bayesian methods offer a formal framework to quantify clinical opinions and use them
308 when relevant historical data is not available. It also provides an opportunity to assess how
309 using the expert opinions under different assumptions of the community of priors impacts
310 the results. The sensitivity of results to alternative expert opinion derivations were also
311 assessed and were helpful in identifying whether various mean and variance parameter
312 derivation choices impact the posterior estimates in the context of this trial.

313 The primary Bayesian analysis results using the interaction prior suggests there is
314 insufficient evidence in favour of VRT being better than usual care. In this sense, the
315 Bayesian analysis is consistent with the frequentist analysis and suggests further research is
316 needed to claim any subgroup effects. The elicitation results indicated that experts believed
317 VRT was better than usual care, and that the effect would be greatest for patients with
318 more severe comorbidities. These opinions were not supported by the trial data, which
319 suggested VRT was worse than usual care for participants with more severe comorbidities. It
320 is unclear whether this inconsistency is due to 'inaccurate' expert opinions or spurious
321 findings from the frequentist analysis which is susceptible to outliers given the small sample
322 size.

323 Trial reporting guidelines advocate against presenting subgroups estimates if the interaction
324 test is not significant [25]. However, this recommendation also prevents us from
325 understanding the effect of an intervention in subgroups of interest. Using an informative

326 prior for the interaction parameters (i.e. the interaction and interaction-variance priors) and
327 an uninformative prior for the overall mean treatment effect in the Bayesian analysis allows
328 us to draw inferences for both subgroups. By ‘borrowing’ information from the treatment
329 effect in the complementary subgroup, subgroup posterior estimates using the interaction
330 priors were more stable (less extreme point estimates and higher precision) than frequentist
331 subgroup estimates. Our proposed Bayesian subgroup model shifts point estimates and
332 their associated credible intervals towards the overall mean effect, whereas classical
333 frequentist approaches keep the point estimate fixed and adjust for multiple comparisons
334 by making the confidence intervals wider. Analyses modelling the treatment effects in a
335 joint model have been reported to be sensible approaches to multiplicity as multiple
336 treatment inferences are directly incorporated in the model [26, 27]. To maximise the use of
337 the trial data, we encourage methodologists to prospectively consider Bayesian subgroup
338 analyses using expert opinions when several low-powered subgroup analyses are planned.

339 The approach to Bayesian subgroup analysis reported in this paper was similar to that
340 conducted by White et al. [18]. However, and as expected in their results, the trial data
341 dominated the prior as each subgroup effect was individually investigated in an adequately
342 powered trial. Therefore, the expert opinions had very little impact on the posterior results.

343 It remains unclear to what extent a Bayesian subgroup analysis provided an increased
344 power to detect subgroup treatment effects at the cost of bias: introducing prior expert
345 opinions in several different specifications generated more precise posterior estimates and
346 a shift compared to the MLEs of the VerDiCT trial. Furthermore, we have identified experts
347 within the centre of the trial, but this may have introduced bias insofar that our experts
348 probably had similar opinions that do not reflect all experts’ opinions on the effect of VRT.

349 An extension of the proposed methodology could investigate how clinical opinions influence
350 the posterior results of a Bayesian subgroup analysis when the true treatment effect is
351 known. A simulation study could investigate the impact of the interaction effect magnitude
352 and different trial sample sizes on inferences, and furthermore evaluate the statistical
353 properties of the Bayesian subgroup analysis by assessing the trade-off between bias and
354 power to detect a treatment effect.

355 **Conclusions**

356 This Bayesian subgroup approach proves its value in cases where fully powered subgroup analyses
357 are not feasible, time-consuming or too expensive. With limited resources, experts can be elicited,
358 and their opinions can help maximise the trial data. Adding experts' opinions to the analyses could
359 increase the precision of the treatment estimate as we have noted in our motivating case study
360 which could in return increase power to detect an effect. Still, and as any subgroup analysis,
361 researchers need to be cautious in appraising each posterior distribution results, acknowledge the
362 limitations of the findings, and provide supporting or contradictory data from other studies when
363 available.

364 **List of abbreviations**

CABG	Coronary artery bypass graft
HCRW	Health and Care Research Wales
HR	Hazard ratio
HRA	Health Research Authority
MLE	Maximum likelihood estimate
NHS	National Health Service
NIHR	National Institute for Health and Research
PATS	Patient Analysis and Tracking System
UK	United Kingdom
VRT	Volume replacement therapy

365

366 **Declarations**

367 **Ethics approval and consent to participate**

368 This study was approved by the HRA and Health and Care Research Wales
369 (HCRW, 19/HRA/0504) ethics committee. Informed written consent was obtained from all
370 participants for the VeRDICT study. The VeRDICT study was registered online at
371 www.isrctn.com (ISRCTN02159606).

372 **Consent for publication**

373 *Not applicable*

374 **Availability of data and materials**

375 The datasets used during the current study are available from the corresponding author on
376 reasonable request.

377 **Competing interests**

378 The authors declare that they have no competing interests.

379 **Availability of data and materials**

380 The codes, protocols, and datasets used during the current study are available from the
381 corresponding author on reasonable request.

382 **Funding**

383 RT was funded by the National Institute for Health Research Methods Fellowship (NIHR-RM-FI-2017
384 08-017). JMB is an NIHR Senior Investigator to which this fellowship was linked. CR and JMB are part
385 funded by the NIHR Bristol Biomedical Research centre. The VeRDICT trial, designed to examine the

386 effects of VRT, was funded by a grant from the Garfield Weston Trust (Ref. PMS/MMS 07/08 3001)
387 to RA. The views and opinions expressed therein are those of the authors and do not necessarily
388 reflect those of the NIHR, UK NHS or Department of Health and Social Care.

389 Authors' contributions

390 RA lead the VerDiCT trial. RT and CR conceptualised this Bayesian study. RT performed the statistical
391 analyses. RT produced the first draft of the manuscript. RA, JMB and CR provided feedback on drafts
392 of the manuscript. All authors have read and approved the manuscript.

393 Acknowledgements

394 The authors would like to thank the experts for completing the elicitation questionnaire, Ben
395 Gibbison for his help piloting the questionnaire and Daniel Fudulu for facilitating the elicitation.

396

397 References

- 398 1. Fan, J., F. Song, and M.O. Bachmann, *Justification and reporting of subgroup analyses*
399 *were lacking or inadequate in randomized controlled trials*. Journal of Clinical
400 Epidemiology, 2019. **108**: p. 17-25.
- 401 2. Kasenda, B., et al., *Subgroup analyses in randomised controlled trials: cohort study on trial*
402 *protocols and journal publications*. BMJ (Clinical research ed.), 2014. **349**: p. g4539-g4539.
- 403 3. Wang, R., et al., *Statistics in Medicine — Reporting of Subgroup Analyses in Clinical Trials*.
404 New England Journal of Medicine, 2007. **357**(21): p. 2189-2194.
- 405 4. Marshall, S.W., *Power for tests of interaction: effect of raising the Type I error rate*.
406 Epidemiologic perspectives & innovations : EP+I, 2007. **4**: p. 4-4.
- 407 5. O'Hagan, A., *Expert knowledge elicitation: subjective but scientific*. The American Statistician,
408 2019. **73**(sup1): p. 69-81.

- 409 6. Morgan, M.G., *Use (and abuse) of expert elicitation in support of decision making for public*
410 *policy*. Proceedings of the National academy of Sciences, 2014. **111**(20): p. 7176-7184.
- 411 7. Chaloner, K., *Elicitation of prior distributions*. Bayesian biostatistics, 1996: p. 141-156.
- 412 8. O'Hagan, T., *Elicitation*. Significance, 2005. **2**(2): p. 84-86.
- 413 9. Hiance, A., S. Chevret, and V. Levy, *A practical approach for eliciting expert prior beliefs*
414 *about cancer survival in phase III randomized trial*. J Clin Epidemiol, 2009. **62**(4): p. 431-
415 437.e2.
- 416 10. Mason, A.J., et al., *Development of a practical approach to expert elicitation for randomised*
417 *controlled trials with missing health outcomes: Application to the IMPROVE trial*. Clin Trials,
418 2017. **14**(4): p. 357-367.
- 419 11. Hemming, K., et al., *Bayesian Cohort and Cross-Sectional Analyses of the PINCER Trial: A*
420 *Pharmacist-Led Intervention to Reduce Medication Errors in Primary Care*. PLOS ONE, 2012.
421 **7**(6): p. e38306.
- 422 12. Sun, C.Q., et al., *Expert prior elicitation and Bayesian analysis of the Mycotic Ulcer Treatment*
423 *Trial I*. Invest Ophthalmol Vis Sci, 2013. **54**(6): p. 4167-73.
- 424 13. Spiegelhalter, D.J., L.S. Freedman, and M.K. Parmar, *Bayesian approaches to randomized*
425 *trials*. Journal of the Royal Statistical Society: Series A (Statistics in Society), 1994. **157**(3): p.
426 357-387.
- 427 14. Cooke, R., *Experts in uncertainty: opinion and subjective probability in science*. 1991: Oxford
428 University Press on Demand.
- 429 15. Cooke, R.M. and L.L.H.J. Goossens, *TU Delft expert judgment data base*. Reliability
430 Engineering & System Safety, 2008. **93**(5): p. 657-674.
- 431 16. Oakley, J., *SHELF: Tools to Support the Sheffield Elicitation Framework*. R package version
432 1.6.0, 2017.
- 433 17. Genest, C. and J.V. Zidek, *Combining probability distributions: A critique and an annotated*
434 *bibliography*. Statistical Science, 1986. **1**(1): p. 114-135.

- 435 18. White, I.R., S.J. Pocock, and D. Wang, *Eliciting and using expert opinions about influence of*
436 *patient characteristics on treatment effects: a Bayesian analysis of the CHARM trials*. *Stat*
437 *Med*, 2005. **24**(24): p. 3805-21.
- 438 19. Sarkar, K., et al., *Preoperative VolumE Replacement therapy in Diabetic patients undergoing*
439 *coronary artery bypass grafting surgery: results from an open parallel group randomized*
440 *Controlled Trial (VeRDICT)*. *Interact Cardiovasc Thorac Surg*, 2019.
- 441 20. R Core Team, *A language and environment for statistical computing*. *R Foundation for*
442 *Statistical Computing, Vienna, Austria*. 2015.
- 443 21. Carpenter, B., et al., *Stan: A probabilistic programming language*. *Journal of statistical*
444 *software*, 2017. **76**(1).
- 445 22. Johnson, S.R., et al., *Methods to elicit beliefs for Bayesian priors: a systematic review*. *Journal*
446 *of clinical epidemiology*, 2010. **63**(4): p. 355-369.
- 447 23. Aspinall, W. and R. Cooke, *Quantifying scientific uncertainty from expert judgement*
448 *elicitation, in. Risk and uncertainty assessment for natural hazards*, 2013: p. 64.
- 449 24. Oakley, J., A. Daneshkhah, and A. O'Hagan, *Nonparametric prior elicitation using the*
450 *Roulette method*. *School of Mathematics and Statistics, University of Sheffield, UK*, 2010.
- 451 25. European Medicines Agency, *International Conference on Harmonisation of Technical*
452 *Requirements for Registration of Pharmaceuticals for Human Use. ICH Harmonised Tripartite*
453 *Guideline: Statistical Principles for Clinical Trials E9*. 1998.
- 454 26. Gelman, A., J. Hill, and M. Yajima, *Why we (usually) don't have to worry about multiple*
455 *comparisons*. *Journal of Research on Educational Effectiveness*, 2012. **5**(2): p. 189-211.
- 456 27. Sjölander, A. and S. Vansteelandt, *Frequentist versus Bayesian approaches to multiple*
457 *testing*. *European Journal of Epidemiology*, 2019. **34**(9): p. 809-821.

458

459