1    Effect of using the same vs different order for second readings of screening mammograms on rates
2    of breast cancer detection: A randomized clinical trial
3

4    Sian Taylor-Phillips[1]* PhD                s.taylor-phillips@warwick.ac.uk
5    Matthew G Wallis[2] FRCR                matthew.wallis@addenbrookes.nhs.uk
6    David Jenkinson[1] PhD                D.J.Jenkinson@warwick.ac.uk
7    Victor Adekanmbi[1] PhD                V.Adekanmbi@warwick.ac.uk
8    Helen Parsons[1] PhD                H.Parsons@warwick.ac.uk
9    Janet Dunn[1] PhD                J.A.Dunn@warwick.ac.uk
10   Nigel Stallard[1] PhD                N.Stallard@warwick.ac.uk
11   Ala Szczepura[1] DPhil                ala.szczepura@coventry.ac.uk
12   Simon Gates[1] PhD                S.Gates@warwick.ac.uk
13   Olive Kearins[3] MSc                Olive.Kearins@phe.gov.uk
14   Alison Duncan[4] FRCR                alison.duncan@uhcw.nhs.uk
15   Sue Hudson[5] MSc                sue.hudson@pasconsulting.co.uk
16   Aileen Clarke[1] MD                Aileen.Clarke@warwick.ac.uk
17

18   1.    Warwick Medical School, the University of Warwick, Coventry, CV4 7AL, UK

19   2.    Cambridge Breast Unit, Cambridge Universities NHS Foundation Trust and NIHR Cambridge
20   Biomedical Research Centre, Cambridge CB2 0QQ, UK

21   3.    Screening QA Service (Midland & East), Seaton House, City Link, Nottingham, NG2 4LA

22   4.    Warwickshire, Solihull and Coventry Breast Screening Service, University Hospitals Coventry
23   and Warwickshire, Coventry CV2 2DX

24   5.    Peel & Schriek Consulting Limited, London N5 1AU

25   *Corresponding author tel (0044) 7725000262

26   Date of revision: 7th April 2016

27   Word count: 4846

28

29   Effect of using the same vs different order for second readings of screening mammograms on rates
30   of breast cancer detection: A randomized clinical trial

31   **Importance:** Interpreting breast screening mammograms is a difficult repetitive task that can result
32   in missed cancers and false positive recalls. In the UK, two film readers independently evaluate each
33   mammogram to search for signs of cancer, and examine digital mammograms in batches. However,
34   a vigilance decrement (reduced detection rate with time on task) has been observed in similar
35   settings.

36   **Objective:** To determine the effect of changing the order for the second film reader of batches of
37   screening mammograms on rates of breast cancer detection.

38   **Design, Setting, and Participants:** A multi-centre, double-blind, cluster randomised controlled trial
39   conducted at 46 specialised breast screening centres from the National Health Service Breast
40   Screening Programme in England for 1 year, (all between 20th December 2012 and 3rd November
41   2014). 360 readers participated (mean 7.8 readers per centre), 186 radiologists, 143 radiography
42   advanced practitioners and 31 breast clinicians, all fully qualified to report mammograms in the NHS
43   breast screening programme.

44   **Intervention:** The two readers examined each batch of digital mammograms in the same order in
45   the control group and in the opposite order to one another in the intervention group.

46   **Main Outcomes and Measures:** The primary outcome was cancer detection rate; secondary
47   outcomes were rates of recall and disagreements between readers.

48   **Results:** Among 1,194,147 women (mean age 59.3, sd 7.49) who had screening mammograms
49   (596642 in the intervention group; 597505 in the control group), the images were interpreted in
50   37,688 batches [median batch size 35 (IQR 16-46)], with each reader interpreting median 176
51   batches (IQR 96 to 278). After completion of all subsequent diagnostic tests, a total of 10,484 cases
52   of breast cancer were detected (0.88%). There was no significant difference in cancer detection rate
53   [5272 (0.88%) vs 5212 (0.87%), difference 0.011% points 95%CI -0.022 to 0.045], recall rate, [24681
54   (4.14%) vs 24894 (4.17%),  difference -0.030% points 95%CI -0.101 to 0.042) or rate of reader
55   disagreements [20,471 (3.43%) vs 20793 (3.48%), difference -0.048% points 95%CI -0.113 to 0.018,]
56   between intervention and control groups.

57   **Conclusions and Relevance:** Interpretation of batches of mammograms by qualified screening
58   mammography readers using a different order versus the same order for the second reading
59   resulted in no significant difference in rates of detection of breast cancer.

60   **Trial Registration:** www.isrctn.com identifier: ISRCTN46603370

61

**Introduction**

62

63     Breast cancer screening currently detects 8.6 cancers per thousand women screened triennially

64     (equivalent to 18,000 cancers per year) in the UK[1] and 4.2 cancers per thousand women screened

65     annually in the US. [2] However, another 2.9 cancers per thousand women screened in the UK [3]

66     (equivalent to 6,030 cancers per year) [1] and 0.9 cancers per thousand women screened in the US

67     are detected between screening rounds in screened women. [2] These arise through cancers growing

68     between screening rounds, and cancers missed at screening. An additional 3.4% of women in the UK

69     (70,715 each year) [1] and 9.3% of women in the US [2] experience false positive recalls at each

70     screening round.

71     Interpreting screening mammograms is a difficult and repetitive visual search task, where

72     characteristics of cancer are disguised amongst background breast parenchyma resulting in false

73     positive recalls and missed cancers. In similar visual search tasks a vigilance decrement of decreasing

74     detection rates with time on task has been observed in a large number of psychological laboratory

75     experiments [4] [5] for example assembly line inspection tasks, [6] airport baggage screening, [7] driving [8]

76     piloting aeroplanes[9] and operating military drones. [10] An effect similar to the vigilance decrement

77     has been observed when examining tests sets of x-rays including mammograms in laboratory

78     conditions although the phenomenon has not previously been explored in breast screening

79     practice.[11,12]

80     In the UK two film readers independently examine each woman's mammograms for signs of cancer.

81     In this study we investigated whether there is a vigilance decrement to detect cancer in breast

82     screening practice, and whether changing the order in which the two experts examined the batch of

83     mammograms could increase the cancer detection rate, through readers experiencing peak vigilance

84     at differing points within the reading batch when examining different women's mammograms.

85

86 **Methods**

87 *Study Design*

88 The Changing Case Order to Optimise Patterns of Performance in Screening (CO-OPS) pragmatic

89 double blind cluster randomised controlled trial was designed to determine whether there is a

90 vigilance decrement in breast cancer screening, and whether changing the order in which the cases

91 are presented can increase cancer detection rate.

92 Ethical approval was granted by the Coventry and Warwickshire National Health Service (NHS)

93 Research Ethics Committee on 27 June 2012 (Reference 12/WM/0182) and written informed

94 consent sought from the director of breast screening at each centre. The trial protocol is provided in

95 supplement 1 and published elsewhere [13] and statistical analysis plan (supplement 2) was finalised

96 before any data were collected.

97 *Intervention and Outcomes*

98 The study compared two parallel groups, each split into two sub-groups to ensure blinding of the

99 readers. The intervention group involved the two readers reading the batch in the opposite order to

100 each other; one forwards, one in reverse. Hence, the two sub-groups: first reader forwards, second

101 reader reverse, and first reader reverse, second reader forward. The control group required the

102 readers to read the batch in the same order as each other; the sub-groups being either both

103 forwards (which is current practice) or both in reverse (to maintain the blinding of a reader to trial

104 group, as they would be aware that they are reading a batch in reverse). Thus each batch (cluster)

105 was randomised with equal probability to one of four groups.

106 The primary outcome was cancer detection rate, (number of women with cancer detected as a

107 proportion of all women screened) as this is the clinically relevant outcome of interest. Secondary

108 outcomes of recall rate (secondary outcome 1) and rate of disagreement between the readers

109 (secondary outcome 2) are designed to examine the proposed mechanism of action. The idea is that

110    reversing the order for one reader results in high vigilance states occurring for the two readers when

111    examining different women's mammograms, so the cancers are detected by at least one of the

112    experts, as outlined in figure 1. If a reader in a high vigilance state detected a cancer missed by their

113    colleague in a low vigilance state, then this would lead to a disagreement between them. All

114    disagreements are 'arbitrated' either by a third reader or group of readers for the final decision of

115    whether to recall the woman for further tests. Assuming the arbitration process performs better

116    than random chance the increases in disagreements would lead to increases in recall rate and cancer

117    detection rate.

118

119    *Participants*

120    Centres were recruited at radiology meetings, through local radiology, radiography and quality

121    assurance groups, and through direct telephone and email contact. The study comprised 46 breast

122    screening centres using digital mammography, each consisting of groups of between 1 and 3

123    hospitals sharing the same computer system for storing women's health records. Characteristics of

124    breast screening centres in England which took part in the trial in comparison to those which did not

125    is provided in table e1. The trial ran for 1 year at each centre, with individual centres starting the

126    study when local consent and research and development approvals were obtained, (start dates were

127    all between 20th December 2012 and 4th November 2013). One centre completed only 4 months of

128    the study due to local technical and workforce issues.

129    Mammograms from women attending routine breast cancer screening at these centres during the

130    study period were included. These were arranged into batches of around 40 women as is standard

131    practice in the UK, and all mammograms taken during the study period were included in the trial,

132    regardless of when they were examined. Each batch contained all cases from a single mammography

133    acquisition machine in a single day. Informed consent was at the centre level, with consent of

134    individual women considered impractical for this system level intervention. In the UK women age 50

135 to 70 are invited to breast screening every 3 years, this study also includes women aged 47-49 and

136 71-73 participating in the age extension trial (NCT01081288), and a small proportion of older women

137 (2.3% of women in the trial) who self-refer as part of the programme.  Women who presented to

138 clinics symptomatically and for high familial risk were excluded.

139 All readers undergo formal training and are accredited by the NHS Breast Screening Programme.

140 They are required to read a minimum of 5,000 cases per year, participate in assessment clinics,

141 formally audit their own performance against their peers,  and maintain ongoing professional

142 development including participating annually in the Personal Performance in Mammographic

143 Screening (PERFORMS) test set.[14] Each centre annually measures and reports results against targets

144 including recall rate, cancer detection rate, and small cancer detection rate, and continuously audits

145 performance through monthly review of interval cancers diagnosed symptomatically between

146 screening rounds, [14]  and monthly checks of mammography acquisition and display equipment and

147 reading room background light levels.[15] Each woman's mammograms are examined by two readers

148 co-located in the same breast screening centre. Readers are instructed to examine the batches

149 independently, but can access the other reader's decision by opening the patient records. In 16 of

150 the 46 centres workflow systems were designed to blind reader 2 to the decision of reader 1. All

151 centres used arbitration when the two readers disagreed, with 13 centres using a single 3rd reader,

152 and 33 centres using group consensus of 2 or more readers.

153 *Randomisation and Blinding*

154 The randomisation took place immediately prior to opening each batch for examination using the

155 Intersystems Caché $RANDOM function within the computer software that the UK National Breast

156 Screening Service (NBSS) uses to manage the work. After randomisation the software automatically

157 displayed the cases in the chosen order to the first and second reader. Readers were aware of the

158 reading order but were blinded to trial group. The trial statistician and the women screened were

159 also blinded to trial group. The unit of randomisation was a batch of mammograms, whereas the

160    unit of observation was the individual mammogram. Simple randomisation was used without

161    stratification or minimisation due to the large number of clusters randomised.

162    *Data collection*

163    The data were collected via an adaptation to the NBSS computer system, which created new tables

164    within the software to record data items pertaining to the trial. The outcomes for every woman

165    screened (including both readers' decision, time of decision, and results of all follow-up tests

166    including biopsy) were added to NBSS as part of each centre's annual reporting requirements, to

167    reduce missing data. The data was extracted through NBSS from each centre, exporting data in Excel

168    format. The datasets from each centre were merged using Excel and R [v 3.0.3 in RStudio v

169    0.98.501]. Cancer was defined as needle biopsy or surgery positive for ductal carcinoma in situ or

170    invasive cancer. Recall for further tests was taken directly from NBSS, which records this decision to

171    enable the follow-up appointment to be made. Disagreement was defined by examining whether

172    the recommendation of whether to recall differed between the first and second readers.

173    *Sample Size*

174    Prior to the study (year 2011-12) the breast cancer detection rate in the UK was 7.8 per thousand

175    women screened. [16] Three years of observational data on patterns of cancer detection with time on

176    task was extracted from routine records at eight breast screening centres in one English region.  This

177    suggested that the intervention may result in one extra cancer detected per 2000 women screened,

178    an increase to 8.3 per thousand women screened.  To detect such an increase required a sample size

179    of 501,361 women in each group, using a 5% significance level and 80% power. The trial had a

180    cluster design, the unit of randomisation being the batch, so the sample size needed to be inflated

181    by the design effect. The inter-cluster correlation coefficient was estimated to be 0.002, resulting in

182    a design effect of 1.09, assuming an average cluster size of 40. Hence, the total sample size required

183    was 1,093,780, which is equivalent to the annual caseload of 44 centres. There were no interim

184    analyses or stopping rules.

185    *Statistical analysis*

186    We used multivariable multilevel logistic regression to analyse factors associated with breast cancer

187    detection, recall and disagreement rates due to the hierarchical nature of the datasets. Analysis was

188    intention to treat, with those not receiving the intervention as allocated included in the analysis.

189    However, women lost to follow-up, technical recalls (mammograms were of insufficient quality to

190    read), and second screening of the same woman were excluded. A three level multilevel model for

191    woman screened (level 1) nested in a batch (level 2) and within a center (level 3) was specified. Four

192    models were constructed for each of the rates stated above. The first model, a null model without

193    any variable was specified to decompose the amount of variance that existed at each level, the

194    second model included the intervention only, the third model included adjustment for known

195    factors associated with cancer and recall (woman's age and whether she had previously attended

196    screening) while the fourth model added the intervention to the adjusted model. All multilevel

197    modelling was performed using MlwiN 2.35 [17] called from Stata statistical software for Windows

198    version 14 [18] using runmlwin routine. For the multilevel logistic regression models, (iterative

199    generalized least squares; penalised quasi-likelihood) IGLS PQL2 estimation was used. [19] Two-tailed

200    tests were used, with p values <0.05 considered significant. The fixed effects (i.e. measures of

201    association) are presented as adjusted odds ratios with their corresponding 95% confidence intervals

202    (CIs). Measures of random effects included intra-cluster correlation (ICC) and median odds ratio

203    (MOR). [20, Merlo, #9] The ICC was calculated by the linear threshold according to the formula used by

204    Snijders et. al. [21] while MOR is a measure of unexplained cluster heterogeneity. Methods used for

205    calculating MOR have been described elsewhere [20, Merlo, #9,22]. Positive Predictive Value was also

206    calculated in the intervention and control groups as the proportion of recalled cases in which cancer

207    was detected.

208    The same models were constructed for three pre-defined sub-groups: women aged under 53 (in

209    whom the intervention may be more effective due to higher breast density increasing the task

210    difficulty); the first and last 5 cases in each batch (where any difference in vigilance would be at its

211    maximum in the intervention group); and the first batch of the day (to examine whether the

212    effectiveness of the intervention may be masked by examining a number of batches in succession).

213    An exploratory post-hoc sub group analysis of cases which are not in the first batch of the day for

214    either reader used the same model structure (to investigate intervention effectiveness when readers

215    may be fatigued).

216    An exploratory post-hoc analysis to measure whether there is a vigilance decrement of decreasing

217    sensitivity to detect cancer with time spent on task, the position in the batch (i.e. $1^{st}$, $2^{nd}$, $3^{rd}$ …) was

218    added as a variable to the unadjusted and adjusted models of cancer detection outlined above. For

219    this analysis the cancer detection rate outcome was personalised to the individual reader who first

220    examined the case, so the outcome had an additional requirement of being correctly identified by

221    the first reader for recall, as well as having cancer identified on follow-up tests. The same modelling

222    approach was applied to recall rate, to measure any systematic change with time on task. In this

223    case it was the recall rate for the first reader, rather than overall from the process that was analysed.

224    Further exploratory post-hoc analysis was conducted to determine whether the lack of effect of the

225    intervention was associated with reader 2 not being blinded to the decision of reader 1 at some trial

226    centres. Including only the sub-group of centres in which reader 2 was blinded to the decision of

227    reader 1, cancer detection rates and recall rates in the intervention group were calculated, and

228    compared to those in the control group.

229    **Results**

230    **Flow of Women in the CO-OPS Trial**

231    1,207,633 women were included in the trial, see figure 2. There were three causes of loss to follow

232    up: 258 (0.02%) were recalled for further tests from screening but did not attend, 233 (0.02%) had

233    an inconclusive needle biopsy result but refused further tests, and 298 women (0.02%) had missing

234    data in the NBSS system. An additional 12,426 cases (1.03%) were judged of insufficient quality for

235    analysis (Technical recall) by the first reader so were not read within batch and could not be included

236    in the analysis, and 271 (0.02%) cases were excluded because the same woman had already been

237    screened that year and included in the trial. This occurred primarily when women moved house and

238    GP practice and consequently were re-invited more quickly than intended.

239    The intervention and control groups were well matched for baseline characteristics including the age

240    and previous attendance of the women screened and batch length, as detailed in table 1.

241    Mammograms were examined by 360 qualified readers, of which 186 were radiologists, 143 were

242    radiography advanced practitioners and 31 were breast clinicians. The median batch length was 35

243    cases (quartiles 16 and 46). Each reader examined a median of 5640 cases, (IQR 2599 to 8458), in a

244    median of 176 batches (IQR 96 to 278) including cases in both the intervention and control groups.

245    Between 1 and 26 batches were examined by each reader in a single day (median 2 IQR 1 to 4). Each

246    centre examined between 8152 and 72714 cases (median 25540 cases).

247

248    **Outcomes**

249

250    The primary outcome, cancer detection rate, was 0.88% (5272/596642) in the intervention group

251    and 0.87% (5212/597505) in the control group (difference 0.011% points 95%CI -0.022 to 0.045), see

252    table 2. The intervention did not affect cancer detection rate in the unadjusted (OR=1.01 95%CI 0.96

253    to 1.06) or adjusted models (OR=1.01 95%CI 0.97 to 1.06), see table 3 and e2. In the adjusted model,

254    cancer detection rate increased with each increasing year of age (OR=1.052 95%CI 1.048 to 1.055)

255    and was higher in women who had not previously attended screening (OR=1.73 95%CI 1.62 to 1.86).

256     The intervention also had no effect in any of the sub-groups of younger age, first and last 5 cases in

257     the batch, the first batch of the day for both readers, or in batches examined second in the day or

258     later by both readers in either the adjusted or unadjusted models. For batches read first in each

259     workday by both readers cancer detection rate was 0.83% (580/70071, 95%CI 0.76% to 0.89%) in the

260     intervention group and 0.88% (623/70715, 95%CI 0.81% to 0.95%) in the control group (difference -

261     0.053% points 95%CI -0.149 to 0.043). For batches read second or subsequent in each workday by

262     both readers cancer detection rate was 0.85% (2472/289786, 95%CI 0.82% to 0.89%) in the

263     intervention group and 0.85% (2473/290671, 95%CI 0.82% to 0.88%) in the control group (difference

264     0.002% points 95%CI -0.045 to 0.050).

265

266

267    The intervention did not affect either of the secondary outcomes, recall rate or rate of

268    disagreements. The recall rate was 4.14% (24681/596642) in the intervention group and 4.17%

269    (24894/597505) in the control group (difference -0.030% points 95%CI -0.101 to 0.042), see table 2.

270    The rate of disagreement was 3.43% in the intervention group (20471/596294) and 3.48%

271    (20793/597387) in the control group (difference -0.048% points 95%CI -0.113 to 0.018), as detailed

272    in table 2. The intervention had no effect on recall rate in the unadjusted (OR=0.993 CI 0.974 to

273    1.013) or adjusted (OR=0.997 CI 0.978 to 1.016) models, (see supplementary table e3) or on rate of

274    disagreement in the unadjusted (OR=0.994 CI 0.971 to 1.019) or adjusted model (OR=0.997 CI 0.974

275    to 1.020), see supplementary table e4. Recall rate was higher with each year of age of the woman

276    screened (OR 1.008 CI 1.007 to 1.010), and was higher in women who had not previously attended

277    breast screening (OR=2.89 CI 2.82 to 2.97). Rate of disagreement was also higher for women at their

278    first screening appointment (OR=2.17 CI 2.11 to 2.24) but lower with each year of increasing age of

279    the woman screened (OR=0.994 CI 0.992 to 0.996). The positive predictive value (PPV) was 21.4%

280    (95% CI 20.8% to 21.9%) in the intervention group and 20.9% (95% CI 20.4% to 21.4%) in the control

281    group (difference 0.420% points 95%CI -0.299 to 1.139). The intervention had no effect on any of the

282    sub-groups (Younger women, first and last cases in the batch, first batch of the day, and second or

283    subsequent batch of the day) for either the adjusted or unadjusted models for either recall rate or

284    rate of disagreements. For batches read first in each workday by both readers recall rate was 4.02%

285    (2818/70071, 95%CI 3.88% to 4.17%) in the intervention group and 4.11% (2904/70715, 95%CI

286    3.96% to 4.25%) in the control group (difference -0.085% points 95%CI -0.291 to 0.121), and rate of

287    disagreements was 3.61% (2531/70071, 95%CI 3.47% to 3.75%) in the intervention group and 3.75%

288    (2653/70715, 95%CI 3.61% to 3.89%) in the control group (difference -0.140% points 95%CI -0.336 to

289    0.057). For batches read second or subsequent in each workday by both readers recall rate was

290    4.10% (11868/289786, 95%CI 4.02% to 4.17%) in the intervention group and 4.15% (12068/290671,

291    95%CI 4.08% to 4.22%) in the control group (difference -0.056% points 95%CI -0.159 to 0.046), and

292   rate of disagreements was 3.23% (9359/289785, 95%CI 3.17% to 3.29%) in the intervention group

293   and 3.28% (9533/290670, 95%CI 3.22% to 3.35%) in the control group (difference -0.050% points

294   95%CI -0.141 to 0.041).

295   Exploratory post-hoc analysis showed that cancer detection rate for individual readers did not

296   change with time spent on task, as represented by near identical odds of detecting cancer between

297   the first and fortieth case (OR=0.987, 95%CI 0.929 to 1.048). Results were very similar in the model

298   adjusted for the characteristics of the woman screened (OR=0.995 95%CI 0.938 to 1.055),

299   supplementary table e5.

300   Exploratory post-hoc analysis showed that recall rate for individual readers (the proportion of

301   women that one reader determined should be recalled) reduced with time on task. The odds of

302   recall decreased over the course of examining 40 cases (OR= 0.83, 95% CI 0.81 to 0.85). The

303   reduction was similar in the model adjusted for woman's age and previous attendance (OR=0.89

304   95%CI 0.87 to 0.91), see supplementary table e6. The mean change over the course of 40 cases was

305   a reduction in recall rate from 6.4% (position 1) to 4.6% (position 40), with the trend continuing in

306   longer batches, see figure 3.

307   Further exploratory post-hoc analysis indicated that there was also no effect of the intervention

308   when readers were blinded to one another's decision. 366,824 cases were read in the trial at the 16

309   centres which blind reader 2 to reader 1 decision. In those centres the cancer detection rate was

310   0.88% (1603/181482, 95%CI 0.84% to 0.93%) in the intervention group and 0.87% (1611/185342,

311   95%CI 0.83% to 0.91%) in the control group (difference 0.014% points, 95% CI -0.046 to 0.074).

312   Similarly recall rate was 4.23% (7669/181482 95%CI 4.13% to 4.32%) in the intervention group and

313   4.23% (7847/185342, 95%CI 4.14% to 4.33%) in the control group (difference -0.008% points, 95%CI

314   -0.138 to 0.122).

315

316

**Discussion**

318    We examined whether an intervention to change the order in which readers examine breast

319    screening cases could improve cancer detection rate. We randomised 1.2 million women in batches

320    of approximately 35 to intervention or control groups. The intervention did not influence cancer

321    detection rate, recall rate, or rate of disagreement between readers. There was no pattern of

322    decreasing cancer detection rate with time on task as predicted by previous research on vigilance

323    decrements as a psychological phenomenon. Instead there was a gradual decrease in recall rate,

324    with an increase in positive predictive value and a decrease in false positive recall of women with

325    time on task. This may reinforce and explain previous observational research which identifies that

326    recall rate is reduced when grouping women's cases into batches. [23]

327    This randomised controlled trial in 1.2 million women was adequately powered to answer the

328    research questions, with over half of the English breast screening service taking part. Effects were

329    measured in a wide range of hospitals, increasing generalisability. Integration into the existing

330    computer systems and reporting mechanisms resulted in very little loss to follow-up (less than 0.1%).

331    Design of the trial computer system was iterative with high user involvement, which increased

332    practicality and facilitated recruitment.

333    This study has several limitations. First, the main limitation is that reading conditions were not

334    controlled so whilst effectiveness in screening practice was measured, efficacy in ideal conditions

335    was not evaluated. In this large pragmatic trial we aimed to measure the effects of the intervention

336    applied to current clinical practice in the UK, and we did not control for or measure working

337    conditions, some of which may affect whether there is a vigilance decrement.  Second, all readers

338    would have met the minimum NHSBSP standards for reading volume, although we did not specify or

339    measure the length of each readers work week, the proportion of their time spent working in breast

340    screening or reading mammograms, the number of work hours or type of work activities each day,

341    number of breaks taken or self-perceptions of fatigue. Similarly whilst there are programme wide

342    auditing methods for reader performance,[14] there will also be centre level variation in management

343    of individual performance which we did not record. Third, the trial did not attempt to implement

344    blinding of reader 2 to the decision of reader 1 where this was not standard practice, as limiting

345    reader's access to computerised and paper notes was not considered possible without

346    compromising patient safety. Fourth, 13% of women in the intervention group did not receive the

347    intervention as intended. The trial software automatically detected these events, which occurred

348    when readers manually overrode the case order and revisited the same case or used barcodes to

349    identify individual cases. These women were included in the intention to treat analysis.

350

351    The trial results were unexpected, and contradict previous research on the vigilance decrement in

352    other fields.[5] The vigilance decrement phenomenon has been reported in many peer reviewed

353    publications,[5] but was not observed in this large randomised controlled trial. These previous studies

354    were primarily undertaken in psychology laboratories rather than in real life settings. Gur et al. [24]

355    demonstrated that performance in experimental conditions and in clinical practice may be very

356    different, suggesting that there is a very different set of incentives in these two settings for the

357    reader.  Hancock contends that the vigilance decrement is entirely a phenomenon created by the

358    conditions designed to measure it. [25] Another explanation for not observing any vigilance decrement

359    is simply that the sessions were too short, however, batches of 40 cases take 20-30 minutes to

360    examine,[26] and the vigilance decrement is usually complete 25 to 35 minutes into the task.[5] The

361    experienced specialists in this study could be less prone to a vigilance decrement, as was found in

362    experienced CCTV operators reviewing a test film. [27] The vigilance decrement phenomenon may be

363    associated with an increase in recall threshold rather than a reduction in performance [28] and if

364    readers already have a low recall threshold so are recalling cases with minimal indications of cancer

365    on the mammograms this may translate to an increase in specificity with minimal decrease in

366    sensitivity. In addition, we have not yet tested the secondary outcome of interval cancer rate (rate of

367    cancers detected symptomatically between screening rounds). If there was a pattern in number of

368    interval cancers with time on task then this may provide evidence of a vigilance decrement. This will

369    be investigated through future analysis of 3 year follow-up data. However we are unlikely to observe

370    such a pattern, because interval cancer rate is inversely proportional to cancer detection rate and

371    this does not change with time on task, and because all cases recalled by one reader did receive a

372    reference standard of peers (independent examination by another reader followed by examination

373    by a third reader or group of readers) and 60% received follow up tests which included ultrasound

374    and biopsy as appropriate. Furthermore the increase in recall rate at the beginning of the batch is

375    many times larger than the total number of interval cancers at screening. [1]

376    A reduction in recall rate with time spent on task has not previously been observed in breast cancer

377    screening. However, an observational study has indicated that examining batches of women's

378    mammograms in one sitting, rather than one-by-one reduces the overall recall rate with no change

379    in cancer detection rate. [23]

380    The systematic reduction in recall rate with time on task for an individual reader did not translate

381    into differences between the intervention and control group (double reading) in overall recall rate or

382    rate of disagreements between readers. There are several possible explanations. The mechanism of

383    action is dependent upon the increased recall rates acting upon the same cases in the control group

384    and different cases in the intervention group. However, the situation is complex. Different readers

385    have different recall thresholds, and different abilities to detect each type of mammographic

386    abnormality (eg spiculated masses, asymmetries, architectural distortions etc). Furthermore each

387    mammogram has overlapping tissue and many features which may appear suspicious. Therefore for

388    any particular pair, the increase in cases recalled at the beginning of the session may not manifest in

389    recalling the same cases. If this is the case, then the intervention would not affect overall recall rate,

390    but it would affect who is recalled, with more women recalled at the beginning of the batch in the

391    control group, and recalls spread more evenly throughout the batch in the intervention group.

392

393    The implications for practice are two-fold. Firstly the intervention of two readers examining a batch

394    of mammograms in the opposite rather than the same order is not effective in increasing cancer

395    detection rate. We have found no evidence of harms from the intervention; however some

396    participating readers reported that it was more difficult to examine cases in reverse order as they

397    also had to reverse associated paperwork. This result is only generalizable to population screening

398    programmes which use two readers to examine mammograms separately. These include the UK NHS

399    breast screening programmes where double reading of mammograms was recommended and

400    became mandatory following the transition to fully digital mammography,[14] European population

401    screening programmes where double reading is recommended and implemented, [29] and Australia

402    where double reading is considered preferable[30] because it increases sensitivity,[31] but not

403    mandated. In the US the Mammography Quality Standards Act and the FDA do not require double

404    reading of mammograms, the decision is made by professional societies and individual centres, and

405    in practice it rarely happens.

406    Secondly, for individual readers recall rate decreases with time spent on task for up to 60 cases, with

407    no concurrent change in cancer detection rate. Therefore we suggest that examining cases in

408    batches of up to 60 is likely to be beneficial. This result was found across 360 readers, encompassing

409    more than half of the NHS Breast Screening Programme in England. Therefore it is likely to be

410    generalisable to screening in England, and may be generalisable across all breast screening

411    programmes using batch reading. Examining mammograms in batches is now standard practice in

412    high volume population breast screening programmes worldwide, with evidence that batch reading

413    increases specificity.[23] However batch reading is not always used, particularly when case volumes

414    are low, such as in practices serving smaller populations. Batch reading is routine for other imaging

415    studies not involving direct radiologist/patient contact with radiology information systems designed

416    for this practice.

417

## Conclusion

419 Interpretation of batches of mammograms by qualified screening mammography readers using a

420 different order versus the same order for the second reading resulted in no significant difference in

421 rates of detection of breast cancer.

## Acknowledgements

## Contributions

447 All authors contributed to the design of the study, and the write up. STP led the study including
448 obtaining funding, design, data collection, analysis and write up. HP performed the sample size
449 calculations. STP, SH and DJ collected the data. STP, DJ, VA and NS contributed to analysing the data.

450

451

452 **Conflicts of Interest**

459

460

**References**

1.  Wilson R. *Breast Screening Programme England Statistics for 2013-14.* Screening and Immunisations team, Health and Social Care Information Centre; 18th February 2015. ISBN 978-1-78386-319-8.
2.  NCI-funded Breast Cancer Surveillance Consortium (HHSN261201100031C). Downloaded on 2/10/15 from the Breast Cancer Surveillance Consortium website http://breastscreening.cancer.gov/.
3.  Bennett RL, Sellars SJ, Moss SM. Interval cancers in the NHS breast cancer screening programme in England, Wales and Northern Ireland. *British Journal of Cancer.* 2011;104(4):571-577.
4.  Mackworth N. The breakdown of vigilance durning prolonged visual search. *Quarterly Journal of Experimental Psychology.* 1948;1(1):6-21.
5.  See JE, Howe SR, Warm JS, Dember WN. Meta-analysis of the sensitivity decrement in vigilance. *Psychological Bulletin.* 1995;117(2):230.
6.  Badalamente RV, Ayoub MM. A Behavioral Analysis of an Assembly Line Inspection Task. *Human Factors: The Journal of the Human Factors and Ergonomics Society.* 1969;11(4):339-352.
7.  Basner M, Rubinstein J, Fomberstein KM, et al. Effects of Night Work, Sleep Loss and Time on Task on Simulated Threat Detection Performance. *Sleep.* 2008;31(9):1251-1259.
8.  Verster JC, Roth T. Vigilance decrement during the on-the-road driving tests: The importance of time-on-task in psychopharmacological research. *Accident Analysis & Prevention.* 2013;58:244-248.
9.  Wiggins MW. Vigilance decrement during a simulated general aviation flight. *Applied Cognitive Psychology.* 2011;25(2):229-235.
10. Gunn DV, Warm JS, Nelson WT, Bolia RS, Schumsky DA, Corcoran KJ. Target Acquisition With UAVs: Vigilance Displays and Advanced Cuing Interfaces. *Human Factors: The Journal of the Human Factors and Ergonomics Society.* 2005;47(3):488-497.
11. Taylor-Phillips S, Elze MC, Krupinski EA, et al. Retrospective Review of the Drop in Observer Detection Performance Over Time in Lesion-enriched Experimental Studies. *Journal of digital imaging.* 2015;28(1):32-40.
12. Gale A, De Silva E, Walker G, Roebuck E, Worthington B. Vigilance decrement and radiological reporting. *Megaw (Ed) Contemporary Ergonomics.* 1989:461-467.
13. Taylor-Phillips S, Wallis MG, Parsons H, et al. Changing case Order to Optimise patterns of Performance in mammography Screening (CO-OPS): study protocol for a randomized controlled trial. *Trials.* 2014;15:17.
14. Wilson R, Liston J. *Quality Assurance Guidelines for Breast Cancer Screening Radiology: NHS Breast Screening Programme Publication Number 59.* Sheffield: NHS Cancer Screening Programmes;2011.
15. Baxter G, Jones V, Milnes V, et al. *Guidance notes for equipment evaluation and protocol for user evaluation of imaging equipment for mammographic screening and assessment: NHSBSP Equipment Report 1411 version 4.* UK: Public Health England; September 2014 2014.
16. McConnell C. *Breast Screening Programme England Statistics for 2011-12.* Screening and Immunisations team, Health and Social Care Information Centre; 27 February 2013. ISBN 978-1-84636-837-0.
17. Rasbash J, Steele, F, Browne W, Prosser B. *A user's guide to MLwiN.* 2015.
18. *Stata statistical software* [computer program]. College Station, TX: StataCorp.; 2015.
19. Goldstein H. *Multilevel statistical models.* London: Hodder Arnold; 2003.
20. Merlo J, Chaix B, Yang M, Lynch J, Råstam L. A brief conceptual tutorial of multilevel analysis in social epidemiology: linking the statistical concept of clustering to the idea of contextual phenomenon. *Journal of epidemiology and community health.* 2005;59(6):443-449.

511   21.   Snijders T, Bosker R. *Multilevel analysis - an introduction to basic and advanced multilevel*
512         *modelling.* Thousand Oaks, California: SAGE publications; 1999.
513   22.   Larsen K, Merlo J. Appropriate assessment of neighborhood effects on individual health:
514         integrating random and fixed effects in multilevel logistic regression. *American journal of*
515         *epidemiology.* 2005;161(1):81-88.
516   23.   Burnside ES, Park JM, Fine JP, Sisney GA. The Use of Batch Reading to Improve the
517         Performance of Screening Mammography. *American Journal of Roentgenology.*
518         2005;185(3):790-796.
519   24.   Gur D, Bandos AI, Cohen CS, et al. The "Laboratory" Effect: Comparing Radiologists'
520         Performance and Variability during Prospective Clinical and Laboratory Mammography
521         Interpretations. *Radiology.* 2008;249(1):47-53.
522   25.   Hancock P. In search of vigilance: the problem of iatrogenically created psychological
523         phenomena. *The American Psychologist.* 2013;68(2):97-109.
524   26.   Taylor-Phillips S, Wallis MG, Gale AG. Should previous mammograms be digitised in the
525         transition to digital mammography? *European radiology.* 2009;19(8):1890-1896.
526   27.   Donald F, Donald C, Thatcher A. Work exposure and vigilance decrements in closed circuit
527         television surveillance. *Applied Ergonomics.* 2015;47:220-228.
528   28.   Broadbent DE, Gregory M. Effects of noise and of signal rate upon vigilance analysed by
529         means of decision theory. *Human Factors: The Journal of the Human Factors and Ergonomics*
530         *Society.* 1965;7(2):155-162.
531   29.   Perry N, Broeders M, de Wolf C, Törnberg S, Holland R, von Karsa L. European guidelines for
532         quality assurance in breast cancer screening and diagnosis. Fourth edition—summary
533         document. *Annals of Oncology.* 2008;19(4):614-622.
534   30.   *Policy on Mammography Screening for Breast Cancer Version 2.* The Royal Australian and
535         New Zealand College of Radiologists;2014.
536   31.   Taylor P, Potts HWW. Computer aids and human second reading as interventions in
537         screening mammography: Two systematic reviews to compare effects on cancer detection
538         and recall rate. *European Journal of Cancer.* 2008;44(6):798-807.

539
540

541    **Table 1 Baseline characteristics for intervention and control groups**

542

|  | Intervention | Control |
|---|---|---|
| **Individual level** | | |
| Mean age of women screened (sd) | 59.3 (7.48) | 59.3 (7.49) |
| Number who have previously attended screening (%) | 126,490 / 596,642 (21.2%) | 128,217 / 597,505 (21.5%) |
| **Cluster level** | | |
| Median batch length (quartiles) | 35 (16,46) | 35 (16,45) |
| Median number of screenings examined by each reader (quartiles) | 2,848 (1,469 , 4,385) | 2,891 (1,543 , 4,458) |
| Median number of batches examined by each reader (quartiles) | 86 (52,143) | 91 (51,138) |
| Median number of screenings examined at each centre (quartiles) | 12,496 (8,997 , 16,523) | 12,908 (9,529 , 16,418) |
| Median number of batches examined at each centre (quartiles) | 376 (282,502) | 364 (272,521) |

543
544

545    Table 2. Primary and secondary outcomes in intervention and control groups.

546    Cancer detection rate, recall rate and rate of disagreement between readers in screenings of
547    previous attenders, screenings of previous non-attenders and all screenings.

548

| Outcome | Intervention | Control | Difference |
|---|---|---|---|
| **Primary outcome: Cancer detection rate** | | | |
| All screenings (CI) | 0.88% (0.86% - 0.91%) 5,272/596,642 | 0.87% (0.85 - 0.90%) 5,212/597,505 | 0.011% points (-0.022 - +0.045) |
| Screenings of previous attenders (CI) | 0.90% (0.87% - 0.92%) 4,214/470,152 | 0.88% (0.85% - 0.91%) 4,122/469,288 | 0.018% points (-0.020 - +0.056) |
| Screenings of previous non-attenders (CI) | 0.84% (0.79% - 0.89%) 1,058/126,490 | 0.85% (0.80% - 0.90%) 1,090/128,217 | -0.014% points (-0.085 - +0.057) |
| **Secondary outcome: Recall Rate** | | | |
| All screenings (CI) | 4.14% (4.09% - 4.19%) 24,681/596,642 | 4.17% (4.12% - 4.22%) 24,894/597,505 | -0.030% points (-0.101 - +0.042) |
| Screenings of previous attenders (CI) | 3.15% (3.10% - 3.20%) 14,819/470,152 | 3.17% (3.12% - 3.22%) 14,869/469,288 | -0.016% points (-0.087 - +0.054) |
| Screenings of previous non-attenders (CI) | 7.80% (7.65% - 7.94%) 9,862/126,490 | 7.82% (7.67% - 7.97%) 10,025/128,217 | -0.022% points (-0.231 - +0.186) |
| **Secondary outcome: Disagreement rate between readers** | | | |
| All screenings (CI) | 3.43% (3.39% - 3.48%) 20,471/596,294 | 3.48% (3.43% - 3.53%) 20,793/597,387 | -0.048% points (-0.113 - +0.018) |
| Screenings of previous attenders (CI) | 2.73% (2.69% - 2.78%) 12,850/469,869 | 2.76% (2.71% - 2.80%) 12,937/469,215 | -0.022% points (-0.088 - +0.044) |
| Screenings of previous non-attenders (CI) | 6.03% (5.90% - 6.16%) 7,621/126,425 | 6.13% (6.00% - 6.26%) 7,856/128,172 | -0.101% points (-0.287 - +0.084) |

549

23

550 **Table 3: Factors associated with cancer detection rate identified by multilevel logistic regression**
551 **models, unadjusted and adjusted for age and previous attendance.**

| Variable | Unadjusted Model OR (CI) | Adjusted model OR (CI) |
|---|---|---|
| **FIXED-EFFECTS (measures of association)** | | |
| **Treatment variable** | | |
| Treatment (vs. control) | 1.01(0.96-1.06) | 1.01(0.97-1.06) |
| **Background factors** | | |
| Age (per year of age) | | 1.052(1.048-1.055) |
| No previous attendance | | 1.73(1.62-1.86) |
| **RANDOM-EFFECTS (measures of variation)** | | |
| **Centre level** | | |
| Variance (SE) | 0.058(0.012-0.104) | 0.038(0.011-0.064) |
| Intra-centre correlation (%) | 1.39 | 0.96 |
| MOR | 1.26 | 1.20 |
| Wald statistics (p-value) | 0.014 | 0.006 |
| **Batch level** | | |
| Variance (SE) | 0.809(0.754-0.863) | 0.595(0.543-0.647) |
| Intra-batch correlation (%) | 20.85 | 16.13 |
| MOR | 2.35 | 2.08 |
| Wald statistics (p-value) | <0.001 | <0.001 |

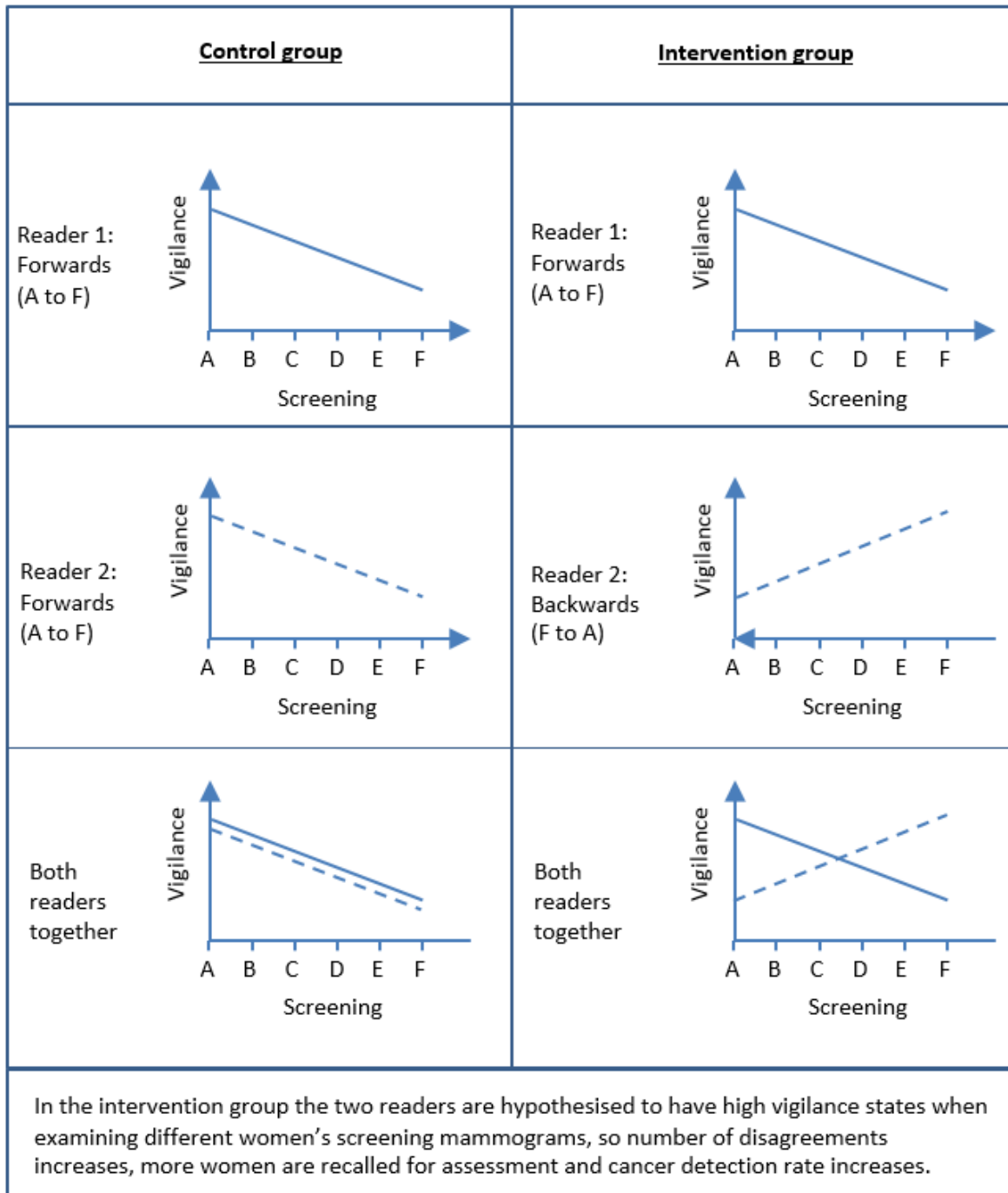552 Abbreviations: SE; standard error, CI; confidence interval, MOR; median odds ratio
553

554

555

**Figure Legends**

557 Figure 1: Proposed mechanism of action of changing case order intervention, assuming the

558 hypothesised vigilance decrement. Each screening represents examining a set of four mammograms,

559 mediolateral oblique and craniocaudal views of both breasts for one woman.

560



561

562

563

564    Figure 2: Study flow of trial comparing same vs different order for presenting batches of
565    mammograms to breast screening readers.

566         [a] Each screening included 4 mammograms (mediolateral oblique and craniocaudal views of
567         both breasts).

568         [b]For each screening there may be multiple reasons why they did not receive the allocated
569         intervention.

570         [c]National Breast Screening Service (NBSS) records are the electronic health records of
571         women screened.

572

80 Breast Screening Centres in England

**Enrollment of centres**

46 centres participated
6 declined (4 were not interested, 2 were already using the intervention)
6 did not meet eligibility criteria
22 were not successfully contacted

Randomized
37,724 batches (n=1,207,633 screenings[a] batch size: mean 32.0, median 35, range 1 to 111)

**Allocation**

Allocated to Intervention Group 18,797 batches (n=603,528 screenings, batch size: mean 32.1, median 35, range 1 to 107)

Received allocated intervention (n=523,781)
Did not receive allocated intervention[b] but included in analysis (n=79,747 screenings):
    Not read in intended order (n=51,599 screenings)
    Reader trainee (n=26,110 screenings)
    Results entered by administrator not reader (n=1 screening)
    Only one reader (n=895 screenings)
    Read using bar code not ordered list (n=16,952 screenings)
    No readers (n=1 screening)

Allocated to Control Group 18,927 batches (n=604,105 screenings, batch size: mean 31.9, median 35, range 1 to 111).

Received allocated intervention (n=559,004)
Did not receive allocated intervention[b] but included in analysis (n= 45,101 screenings):
    Not read in intended order (n=40,528 screenings)
    Results entered by administrator not reader (n=1 screenings)
    Only one reader (n=625 screenings)
    Read using bar code not ordered list (n=17,176 screenings)
    No readers (n=0 screenings)

**Follow-Up**

Lost to follow-up:
    NBSS records[c] not updated (n=172 screenings) Inconclusive biopsy and did not attend any further test (n=115 screenings)
    Recalled for assessment but did not attend appointment (n=118 screenings)
    Discontinued intervention (n=0 screenings)

Lost to follow-up:
    NBSS records[c] not updated (n=126 screenings)
    Inconclusive biopsy and did not attend any further test (n=118 screenings)
    Recalled for assessment but did not attend appointment (n=140 screenings)
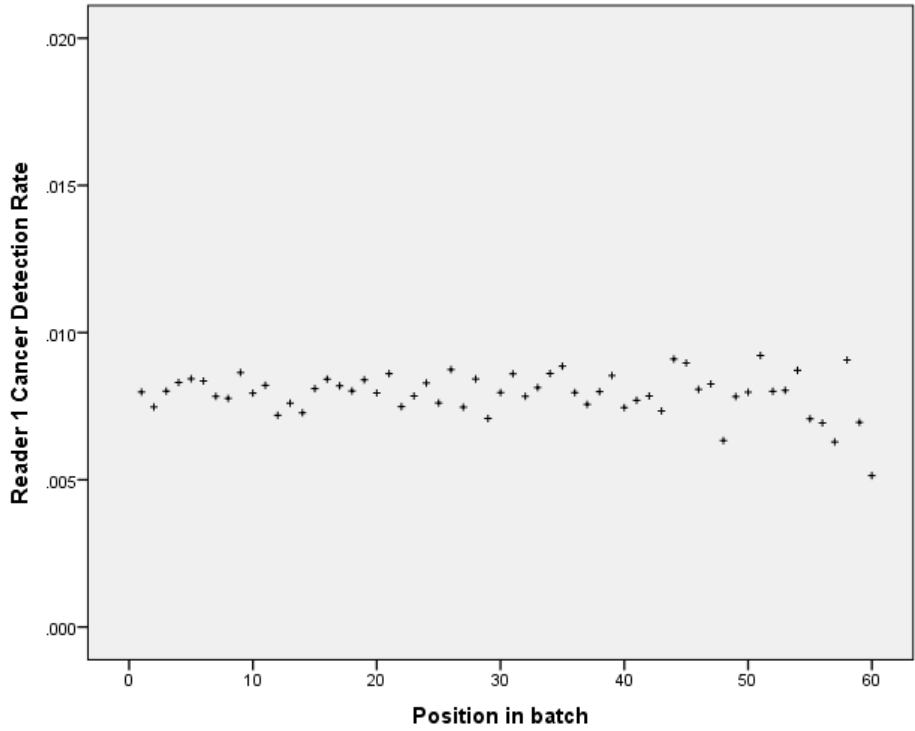    Discontinued intervention (n=0 screenings)

**Analysis**

Analysed 596,642 screenings in 18,779 batches (batch size: mean 31.8, median 35, range 1 to 106)
Excluded from analysis:
    Technical Recall (n=6,339 screenings)
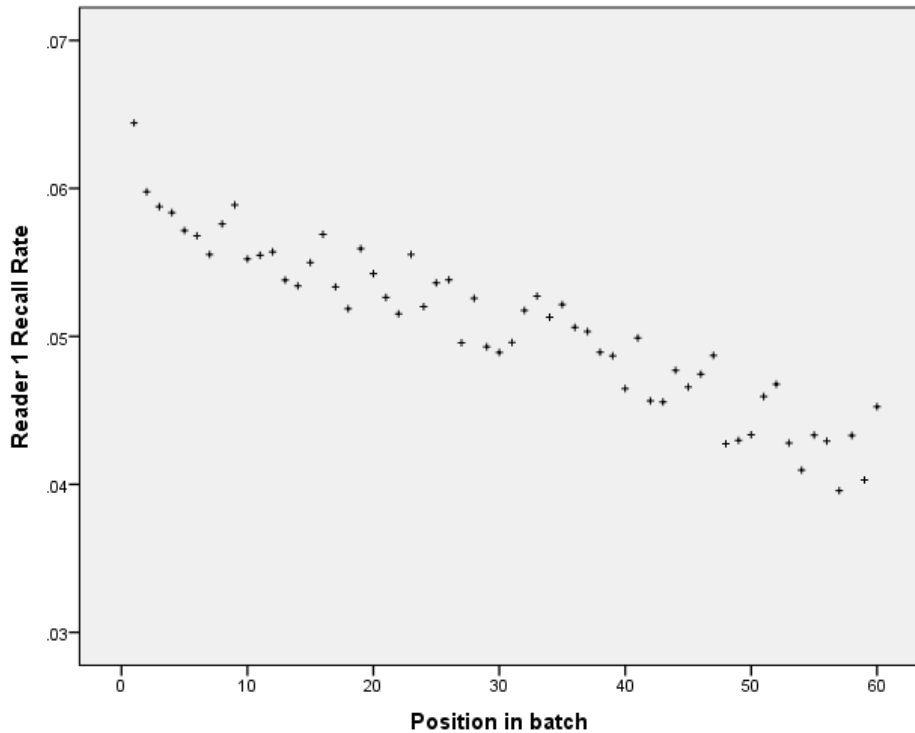    Subsequent screen of same woman (n=142)

Analysed 597,505 screenings in 18,909 batches (batch size: mean 31.6, median 34, range 1 to 110)
Excluded from analysis:
    Technical Recall (n=6,087 screenings)
    Subsequent screen of same woman (n=129)

577    Figure 3 Average patterns of cancer detection rate and recall rate for a single reader over the course
578    of examining a batch of mammograms.

579    Each data point represents the mean recall or cancer detection rate over all cases examined by
580    reader 1 at that position in the batch. 1,173,930 cases are included, examined as reader 1 by 348
581    readers, median number of screenings per batch position is 21,931 (IQR 10,133 to 28,126).

582



583
584
585

28