

Optimizing the pairs of radiologists that double read screening mammograms

Radiology 2023; 309(1):e222691

<https://doi.org/10.1148/radiol.222691>

Jessie J.J. Gommers, MSc; Craig K. Abbey, PhD; Fredrik Strand, MD, PhD; Sian Taylor-Phillips, PhD; David J. Jenkinson, PhD; Marthe Larsen, MSc; Solveig Hofvind, PhD; Ioannis Sechopoulos, PhD; Mireille J.M. Broeders, PhD

From the department of Medical Imaging, Radboud University Medical Center, Geert Grooteplein 10, 6525 GA, Post 766, Nijmegen, The Netherlands (J.J.J.G., I.S.); Department of psychological and brain sciences, University of California, Santa Barbara, United States (C.K.A.); Department of Oncology-Pathology, Karolinska Institute, Stockholm, Sweden (F.S.); and Breast Radiology, Karolinska University Hospital, Stockholm, Sweden (F.S.); Warwick Medical School, University of Warwick, Coventry, United Kingdom (S.T.P., D.J.J.); Section for Breast Cancer Screening, Cancer Registry of Norway, Oslo, Norway (M.L., S.H.); Department of Health and Care Sciences, UiT The Arctic University of Norway, Tromsø, Norway (S.H.); Dutch expert center for screening (LRCB), Nijmegen, The Netherlands (I.S., M.J.M.B.); Technical Medicine Center, University of Twente, Enschede, The Netherlands (I.S.); Department for Health Evidence, Radboud University Medical Center, Nijmegen, The Netherlands (M.J.M.B.)

Address correspondence to M.J.M.B.

+316 15 38 50 63

Mireille.Broeders@radboudumc.nl

Department for Health Evidence, Radboud University Medical Center, Geert Grooteplein 10, 6525 GA, Post 766, Nijmegen, The Netherlands

Funding: aiREAD – Accurate and Intelligent Reading for EARlier breast cancer Detection (project number 17912) supported by the Dutch Research Council (NWO), the Dutch Cancer Society (KWF), and Health Holland (HH).

Manuscript Type: Original Research

Data sharing statement: Data analyzed during the study were provided by third parties. Requests for data should be directed to the providers indicated in the Acknowledgements.

Summary statement:

Performance characteristics of mammography readers influenced the performance of pairs, but specific pairing strategies did not result in significantly different overall performance compared to that resulting from random pairing strategies.

Key results:

- Retrospective data (3,592,414 exams) from three population-based breast cancer screening programs (Sweden, England, and Norway) showed variation in cancer detection and abnormal interpretation rates among radiologists.
- Performance of specific pairs of radiologists was influenced by what types of individual readers were involved.
- Specific radiologist pairing strategies were not significantly different from the random radiologist pairing strategies. Data in which all radiologists read all examinations is needed to explore if there is an optimal pairing strategy that maximizes performance.

Abbreviations:

AIR = abnormal interpretation rate

CDR = cancer detection rate

FN = false negative

FP = false positive

HH = reader characterized by a high cancer detection rate & high abnormal interpretation rate

HL = reader characterized by a high cancer detection rate & low abnormal interpretation rate

LH = reader characterized by a low cancer detection rate & high abnormal interpretation rate

LL = reader characterized by a low cancer detection rate & low abnormal interpretation rate

TN = true negative

TP = true positive

Abstract

Background:

Despite variation in performance characteristics among radiologists, the pairing of radiologists for the double reading of screening mammography is performed randomly. It is unknown how to optimize pairing to improve screening performance.

Purpose:

To investigate whether radiologist performance characteristics can be used to determine the optimal set of pairs of radiologists to double read screening mammograms for improved accuracy.

Materials and Methods:

This retrospective study was performed with reading outcomes from breast cancer screening programs in Sweden (2008-2015), England (2012-2014), and Norway (2004-2018). Cancer detection rates (CDR) and abnormal interpretation rates (AIR) were calculated, with AIR determined by either reader flagging a case as abnormal. Individual readers were divided into performance categories based on their high/low CDR and AIR. The performance of individuals determined the classification of pairs. Random pair performance, for which any type of pair was equally represented, was compared to the performance of specific pairing strategies, which consisted of pairs with readers that are either opposite or similar in AIR and/or CDR.

Results:

Based on a minimum number of examinations per reader/pair, the final study sample consisted of 3,592,414 examinations (Sweden: n:965,263; England:837,048; Norway:1,790,103). The overall AIRs and CDRs for all specific pairing strategies (Sweden: AIR range: 45.5 – 56.9/1,000, CDR range: 3.1 – 3.6/1,000, England: AIR range: 68.2 – 70.5/1,000, CDR range: 8.9 – 9.4/1,000, Norway: AIR range: 81.6 – 88.1/1,000, CDR range: 6.1 – 6.8/1,000) were not significantly different from the random pairing strategy (Sweden: AIR: 54.1/1,000, CDR: 3.3/1,000, England: AIR: 69.3/1,000, CDR: 9.1/1,000, Norway: AIR: 84.1/1,000, CDR: 6.3/1,000).

Conclusion:

Pairing a set of readers based on different pairing strategies did not show a significant difference in screening performance when compared to random pairing. Future studies should include datasets with more than two readers reading the examinations to explore the possibility of improving screening performance by pairing.

Introduction

Population-based breast cancer screening programs with mammography have proven to be effective in reducing breast cancer-specific mortality (1, 2). Nevertheless, breast cancer is the most commonly diagnosed cancer and a leading cause of cancer death among women worldwide (3). Radiologists miss 3%–40% of the mammographically visible cancers (4-6). At the same time, false positive screening results lead to unnecessary workup, participant anxiety, and a reduction in cost-effectiveness (7).

A potential avenue for reducing the rate of errors in a screening program may be to optimize the double reading of screening mammograms. Double reading facilitates the interpretation of mammograms by two individuals with different cognitive, perceptual, and decision-making expertise. Previous studies have shown that double reading increases the cancer detection rate when compared to single reading (8-10). As a result, the European Commission Initiative on Breast Cancer (ECIBC) recommends double reading and breast cancer screening programs in Australia, Europe, and New Zealand have implemented double reading (11). For breast cancer screening programs where mammograms are currently single read, it is possible that artificial intelligence (AI) may be added as a second reader in the future (12-15).

The pairing of radiologists that double read screening mammograms is currently assigned randomly, out of convenience, or to balance the workload. However, screening performance among radiologists varies (16). A previous multi-reader multi-case study with an enriched case set in a laboratory setting demonstrated that it was possible to improve the accuracy of mammography interpretation with double reading when the set of paired radiologists was optimized (17). Optimal pairing may thus be feasible, but no data were available on what factors determined the optimized set of pairs. To our knowledge, only one previous study investigated what prospective criteria could be used to pair radiologists optimally (18). In that study, Gandomkar et al. suggested to pair radiologists with different cognitive eye-tracking metrics to optimize the pairings. However, eye-tracking data is not yet routinely available in screening programs.

Therefore, in our study, we attempted to determine if pairing optimization can be achieved based on the individual readers' performance characteristics that are routinely available in the screening program. The optimal set of pairs is defined as the one that results in the best overall screening performance, characterized by a high cancer detection rate and a low abnormal interpretation rate, for the entire case set. However, it is not currently known how to prospectively select the optimal set of pairs of radiologists. The pair composed of the two most accurate radiologists of the whole program may yield the best overall outcome if they read all examinations, but this is not realistic since the case load needs to be divided evenly. In this study we aim to investigate whether radiologist performance characteristics can be used to determine the optimal set of pairs within a group of radiologists to double read women's screening mammograms.

Materials and Methods

Our retrospective study was performed with de-identified, retrospectively collected screening reading outcomes. We used three datasets: the Swedish CSAW (Cohort of Screen-Age Women) dataset (19), the CO-OPS (Changing case Order to Optimize patterns of Performance in Screening) dataset from England (20, 21), and registry data from BreastScreen Norway to be able to compare

different screening practices. All analysis was performed separately on each dataset, with no pooling done at any point. The CSAW and CO-OPS datasets were used in previously published works (15, 19-26), but these did not investigate different pairing strategies. The use of the CSAW dataset for the purpose of research has been approved by the regional Ethical Review Board, which waived the requirement for written informed consent. For the CO-OPS dataset, ethical approval for the original trial was obtained from Coventry and Warwickshire National Health Service Research Ethics Committee and informed consent was obtained from each director of breast screening. The Norwegian dataset was disclosed with legal bases in the Cancer Registry of Norway Regulations of 21 December 2001 No. 47 and no approval by an ethical board or informed consent was needed as the project was considered a research quality assurance project.

Screening procedures

The screening programs of the three countries differ in several aspects (Figure 1).

The CSAW dataset consists of women from Stockholm county who attended mammography screening between 2008 and 2015. Details of the dataset have been described elsewhere (19, 22). Women 40 to 74 years of age were invited for two-view digital screening mammography every 18 or 24 months. Women >49 years old were invited every 24 months, whereas younger women were invited every 18 months. The examinations were assessed by independent double reading. In most centers the second radiologist was blinded to the assessment performed by the first radiologist. Screening examinations with concordant negative assessments were considered normal. Examinations with one or two abnormal assessment(s), were flagged for consensus, where the final recall decision was made.

The CO-OPS dataset includes women (predominantly aged 47-73 years) invited to two-view digital mammography screening between 2012 and 2014 at 46 breast screening centers throughout England. Women were invited every three years and the mammograms were assessed by two expert readers (radiologists, advanced radiography practitioners, or breast clinicians), who independently decided whether the woman should be recalled. There are local variations in practice regarding blinding, but at most centers the second reader was not blinded to the decision of the first. Screening examinations with concordant negative assessments were considered normal and examinations with concordant positive readings were recalled. Discrepant readings were resolved through a single third reader or group arbitration. In some centers with high recall rates, arbitration was also applied when both readers indicated recall. The original trial and protocol are published elsewhere (20, 21).

The dataset from BreastScreen Norway includes data from women who attended digital mammography screening between 2004 and 2018. BreastScreen Norway offers women aged 50- to 69 years biennial two-view digital mammography. The screening examinations were independently assessed by two radiologists who were blinded to each other's assessment (27). Both radiologists assigned a score from 1 to 5, indicating the suspiciousness of mammographic findings (1 - negative for malignancy, 5 - high suspicion of malignancy). If both readers assigned a score of 1, the screening was assigned as being normal. If either or both radiologist(s) assigned a score of 2 or higher, a consensus meeting determined whether the woman should be recalled. A score of 2 or higher was also the threshold used in this study as a positive assessment.

Study population

All datasets consisted of more than 1 million screening examinations (Table 1). Screening performance, based on the final recall decision, varied among the datasets. The differences were at least partly to be expected, due to differences in screening policies. The recall- and cancer detection rate of the Swedish dataset were the lowest and the recall- and cancer detection rate of the English dataset were the highest.

Breast cancer was defined as needle biopsy or surgery samples that tested positive for ductal carcinoma *in situ* or invasive cancer and was either diagnosed after further assessment in screening or clinically before the next screening examination (i.e., interval cancer). For the English and Norwegian dataset any breast cancer detected before the next screening examination was used for the analyses. For the Swedish dataset breast cancers detected within 18 and 24 months after screening for women aged ≤ 49 years and >49 years, respectively, were used as we did not have information on the exact date of the next screening examination.

Analyses were performed on each dataset separately. Reliable performance measures in a screening program with a low event rate can only be established in sufficiently large datasets (28). Data from readers with less than 500 interpretations in total, readings with unknown reader data, readings with recall because of clinical symptoms or inadequate images, and pairs with a relatively low volume in the dataset ($<$ mean number of interpretations per pair) were excluded. The mean number of interpretations per pair was used to have one consistent selection criterion for all three datasets.

Statistical analysis

We calculated the cancer detection rate (CDR) (true-positives/1,000 readings) and abnormal interpretation rate (AIR) ((true-positives + false-positives)/1,000 readings) for the individual readers and pairs and hypothesized the performance of different pairing strategies (Figure 2). For the purposes of this study, AIR is referred to as the decision(s) of the individual or the pair of readers, while recall rate refers to the final decision according to the program policy, including consensus or arbitration, if necessary.

Individual readers

The CDR and AIR of the individual readers were calculated. The readers were divided into four performance categories, using the individual weighted mean CDR and AIR as cutoffs (top left Figure 2). These cutoffs were calculated for each of the three country datasets separately and the weights were the number of examinations for each reader. Readers were classified into high CDR & low AIR (HL), high CDR & AIR (HH), low CDR & AIR (LL), or low CDR & high AIR (LH).

Pairs

CDR and AIR were also calculated for the pairs of readers. For paired reader assessments we used the pairing rule shown in Figure 3, where a positive assessment was defined by either or both reader(s) flagging an exam as abnormal (i.e., all concordant positive assessments or discrepant assessments were defined as a positive paired reading). By using this pairing rule, we aimed to optimize the performance in terms of cancer cases being, at least, sent to consensus or arbitration after double reading. Consensus discussion or arbitration itself were not considered as we are interested in optimizing the original paired reading only. Pairs were classified based on the performance of the involved individuals. Based on the four different types of individual readers, 16

different types of pairs exist (top right Figure 2). For each of these 16 types of pairs, the average CDR and AIR were calculated by taking the unweighted mean CDR and AIR of the unique pairs involved in that specific type (e.g. the average AIR of the HL+HL pair type was calculated by averaging over the pairing rule defined AIRs of the pairs that were classified as HL+HL).

Hypothetical pairing strategies

Random hypothetical pair performance was compared to the performance of specific hypothetical pairing strategies. For our analyses we were mainly interested in whether pairing readers with 1) opposite or 2) similar performance characteristics was beneficial. Based on our two performance measures (AIR & CDR), we tested six specific pairing strategies against the random pairing strategy (bottom Figure 2):

- opposite AIR and CDR
- opposite CDR
- opposite AIR
- similar AIR and CDR
- similar CDR
- similar AIR

Random pair performance was estimated by averaging over the CDR and AIR of the 16 types of pairs, assuming each of the 16 types was equally represented (top right Figure 2). Specific pair performance was estimated by averaging over the CDR and AIR of the types of pairs that belong to the specific pairing strategy (bottom Figure 2). For example, a pairing strategy with opposite CDR readers in a pair consists of 8 types of pairs (HL&LL, LL&HL, HL&LH, LH&HL, HH&LH, LH&HH, HH&LL, LL&HH), which all involve one reader characterized as high CDR and one reader characterized as low CDR. The grouped screening performance measures of the six specific pairing strategies were compared to the performance measures of the random pairing strategy.

Bootstrap resampling of the screening examinations with corresponding readers (n=1,000) was used to obtain 95% confidence intervals. For each bootstrap sample all analyses steps were performed as described above, including the assessment of individual and paired performance, the classification of readers and pairs, and the evaluation of the hypothetical pairing strategies. Bonferroni-corrected *P*-values <0.008 (0.05/6) were regarded as statistically significant and all statistical analyses were performed in R studio version 4.1.0 (RStudio, PBC).

Results

The final study sample of the Swedish, English, and Norwegian dataset included n=965,263, n=837,048, and n=1,790,103 screening examinations, respectively (Figure 4).

Study sample characteristics

The study subsamples involved different numbers of readers and pairs with different screening performance (Table 2). The Swedish subsample included the youngest study population. The readers in Sweden had the lowest individual weighted mean AIR and CDR (36.3 and 3.1 per 1,000 respectively). The readers in the English subsample had the highest individual CDR of 8.3 per 1,000.

Paired AIR and CDR were higher than individual AIR and CDR for all three study subsamples. This is explained by the pairing rule for which any discrepant reading was defined as a positive reading, thereby increasing the AIR and CDR of the pairs when there is disagreement. Paired AIR and

CDR were lowest for the Swedish subsample (47.8 and 3.4 per 1,000 respectively), while paired CDR was highest for the English subsample (8.9 per 1,000). The Norwegian subsample showed most disagreement among the readers (5.6%), resulting in the highest paired AIR of 77.0 per 1,000.

Individual performance

Figure 5 shows the classification of the individual readers into four performance groups. The individual weighted mean AIR (Table 2, Sweden: 36.3/1,000, England: 48.2/1,000, Sweden: 49.0/1,000) and CDR (Sweden: 3.1/1,000, England: 8.3/1,000, Norway: 5.2/1,000) were used as cutoffs for the respective datasets.

Paired performance

Figure 6 shows the resulting paired AIRs and CDRs for the sixteen specific pair types. The pattern for all three study subsamples looks similar. Pairs consisting of two high-AIR readers (blue and red) resulted in high paired AIR values (>63, >87, >103 per 1,000 for the Swedish, English, and Norwegian subsamples, respectively), and pairs consisting of two low-AIR readers (green and yellow) resulted in low paired AIR values (<38, <51, <59 per 1,000 for the Swedish, English, and Norwegian subsamples, respectively). The same applies to pairs consisting of two high-CDR readers (green and blue) or two low-CDR readers (yellow and red), resulting in high (>3.8, >10.5, >6.2 per 1,000 for the Swedish, English, and Norwegian subsamples, respectively) or low (<3.1, <7.6, <5.7 per 1,000 for the Swedish, English, and Norwegian subsamples, respectively) paired CDRs, respectively. Pairs consisting of opposite AIR and/or CDR readers resulted in average paired AIRs and CDRs compared to the other pair types.

Hypothetical set of pairs

To find out if individual reader performance characteristics can be leveraged to identify the optimal set of pairs, random hypothetical pair performance was compared to the hypothetical group performance of six specific pairing strategies. The group AIR and CDR of the random pairing strategies were 54.1 per 1,000 (CI: 46.1-62.1) and 3.3 per 1,000 (CI: 2.8-3.9) for the Swedish subsample, 69.3 per 1,000 (CI: 63.7-74.9) and 9.1 per 1,000 (CI: 8.3-9.9) for the English subsample, and 84.1 per 1,000 (CI: 80.3-88.0) and 6.3 per 1,000 (CI: 5.9-6.7) for the Norwegian subsample (Figure 7). The confidence intervals from the grouped AIRs and CDRs of the specific pairing strategies overlapped with those from the random pairing strategy. The group AIRs and CDRs for the specific pairing strategies were thus not statistically significantly different from the random pairing strategy in all three study subsamples (Table 3).

Discussion

Our retrospective study showed that **performance characteristics of mammography readers influenced the performance of pairs, but specific pairing strategies did not result in significantly different overall performance compared to that resulting from random pairing strategies.** The specific pairing strategies included some higher performing pairs as well as lower performing pairs that together balanced the overall screening performance of the pairing strategies to AIR and CDR values that were very similar.

In most countries the Picture Archiving and Communication System provides opportunities for strategic pairing as readers can assess the screening mammograms from different locations. Although a previous study by Brennan and colleagues demonstrated that some pairing schemes

were better than others, their study was not designed to identify the criteria that can be used prospectively to determine the best pairing schemes (17). In our study, we attempted to determine if this pairing optimization can be achieved based on the individual readers' performance, but we were unable to identify a pairing rule that consistently and significantly improved the overall program performance. This could be due to the actual underlying optimal criterion not being the individual performance of the readers, the inconsistencies across screening programs introducing differences in the predicted outcomes, the classification of the readers or the datasets having too few reads per examination to exhaustively explore the different pairing strategies.

In the first place, the true optimal pairing criteria may be related to other factors than individual performance. For example, different readers may be better at detecting specific types of findings (e.g., calcifications vs. soft tissue lesions vs. architectural distortions), and therefore pairing based on those differing abilities would yield the optimal program performance, as opposed to the criteria investigated here. The concept of AI as a second reader is also an upcoming and promising method that has gained significant attention in recent years and this could be a way to incorporate double reading in single-reading breast cancer screening programs. A potential implementation of AI as a second reader could involve adjusting the AI operating point settings to the performance of the paired human reader to counter the reader's operating point and hence optimize screening performance. Future research should therefore focus on what pairing criteria may optimize screening performance for both double reading programs as well as single human reading programs with AI as a second reader.

Secondly, the different policies of the screening programs may introduce differences in performance measures among the three subsamples, which then confound the results of the optimization process investigated here. Swedish readers had the lowest individual AIR and CDR, probably due to the younger group of women who were screened every 18 months. The English readers had the highest CDR, probably because of the screening interval of 36 months. Nevertheless, individual AIR was not the highest for the English subsample, but for the Norwegian subsample. One reason for this may be that radiologists in England are more reluctant to flag a case as abnormal because they know that the woman will be recalled if the other radiologist also decides to recall, whereas in Norway consensus will always decide on recall, even after two positive reader assessments. Paired performance showed that Norwegian pairs disagreed more than Swedish and English pairs. This may be because all Norwegian readers were blinded to each other's assessment, whereas for the Swedish and English subsamples not all readers were blinded.

Furthermore, the performance measures of some of the individual readers are close to the predefined cutoff line. Therefore, although dichotomized as high or low CDR and/or AIR, those readers might actually perform very similarly to others that are just on the other side of the corresponding threshold. Therefore, if there actually were an impact on overall performance by pairing readers with specific CDR/AIR characteristics, these would be challenging to tease out with datasets where the actual CDR/AIR differences are small.

Finally, although this study consisted of large study samples, the individual examinations were read by only one pair of readers, and therefore there might not have been enough pair realizations to identify prospective selection criteria for the optimization of pairs. Therefore, it could be helpful to exhaustively evaluate all theoretically possible pairs. Research should therefore focus

on simulating individual radiologist assessments, making it possible to analyze data consisting of results of all readers reading all cases, and exploring possible optimal pairing strategies.

This study has several strengths and limitations. A major strength is that this study involved data from actual screening programs, in which reading behavior influenced care. Furthermore, this study was able to compare three screening programs. However, our study was affected by incorporation bias, because if a reader selected recall, then cancer was more likely detected because of follow-up tests. This bias was reduced by taking into account interval cancers, and therefore identifying the false negatives in the screening datasets. Whilst this biased overall cancer detection accuracy upwards, it was unlikely to impact the comparisons in this paper as this effect would have applied to all readers. Furthermore, the hypothetical pairing strategies rely on the assumption that each type of reader was equally represented in the pairings. This allowed us to make a fair comparison of the different pairing strategies without being influenced by what type of readers were included (e.g., including more high CDR/low AIR than low CDR/high AIR readers would automatically result in better performance independent of the pairing strategy). In actual screening practice there will be variation in the number and type of readers involved, which should be considered when interpreting our results. In addition, there is some variation in the blinding of the second reader, both across and within the different study subsamples. We did not have information on which readers were blinded and how much experience the readers had, so we could not control for these differences.

In conclusion, this study shows that the type of readers involved in a pair influence paired screening performance. Nevertheless, pairing strategies based on CDR and AIR performance characteristics for the full set of pairs did not consistently show significant differences in paired screening performance. Future studies should include datasets with screening examinations read by more than two readers and test pairing strategies with a variable number of reader types to explore the possibility of improving overall screening performance with different pairing strategies.

Acknowledgements

The CO-OPS dataset was funded by an NIHR Postdoctoral Fellowship and an NIHR Career Development Fellowship (CDF-2016-09-018).

References

1. Broeders M, Moss S, Nyström L, Njor S, Jonsson H, Paap E, et al. The impact of mammographic screening on breast cancer mortality in Europe: a review of observational studies. *J Med Screen*. 2012;19 Suppl 1:14-25.
2. Myers ER, Moorman P, Gierisch JM, Havrilesky LJ, Grimm LJ, Ghatge S, et al. Benefits and Harms of Breast Cancer Screening: A Systematic Review. *Jama*. 2015;314(15):1615-34.
3. Sung H, Ferlay J, Siegel RL, Laversanne M, Soerjomataram I, Jemal A, et al. Global Cancer Statistics 2020: GLOBOCAN Estimates of Incidence and Mortality Worldwide for 36 Cancers in 185 Countries. *CA Cancer J Clin*. 2021;71(3):209-49.
4. Broeders MJ, Onland-Moret NC, Rijken HJ, Hendriks JH, Verbeek AL, Holland R. Use of previous screening mammograms to identify features indicating cases that would have a possible gain in prognosis following earlier detection. *Eur J Cancer*. 2003;39(12):1770-5.
5. Majid AS, de Paredes ES, Doherty RD, Sharma NR, Salvador X. Missed breast carcinoma: pitfalls and pearls. *Radiographics*. 2003;23(4):881-95.

6. Weber RJ, van Bommel RM, Louwman MW, Nederend J, Voogd AC, Jansen FH, et al. Characteristics and prognosis of interval cancers after biennial screen-film or full-field digital screening mammography. *Breast Cancer Res Treat.* 2016;158(3):471-83.
7. Brewer NT, Salz T, Lillie SE. Systematic review: the long-term effects of false-positive mammograms. *Ann Intern Med.* 2007;146(7):502-10.
8. Harvey SC, Geller B, Oppenheimer RG, Pinet M, Riddell L, Garra B. Increase in cancer detection and recall rates with independent double interpretation of screening mammography. *AJR Am J Roentgenol.* 2003;180(5):1461-7.
9. Coolen AMP, Voogd AC, Strobbe LJ, Louwman MWJ, Tjan-Heijnen VCG, Duijm LEM. Impact of the second reader on screening outcome at blinded double reading of digital screening mammograms. *Br J Cancer.* 2018;119(4):503-7.
10. Taylor P, Potts HW. Computer aids and human second reading as interventions in screening mammography: two systematic reviews to compare effects on cancer detection and recall rate. *Eur J Cancer.* 2008;44(6):798-807.
11. ECIBC. Use of double reading in mammography screening 2019 [Available from: <https://healthcare-quality.jrc.ec.europa.eu/ecibc/european-breast-cancer-guidelines/organisation-of-screening-programme/double-reading-in-mammography-screening>].
12. Marinovich ML, Wylie E, Lotter W, Pearce A, Carter SM, Lund H, et al. Artificial intelligence (AI) to enhance breast cancer screening: protocol for population-based cohort study of cancer detection. *BMJ Open.* 2022;12(1):e054005.
13. Sharma N, Ng AY, James JJ, Khara G, Ambrozay E, Austin CC, et al. Retrospective large-scale evaluation of an AI system as an independent reader for double reading in breast cancer screening. *medRxiv.* 2022:2021.02.26.21252537.
14. Larsen M, Aglen CF, Hoff SR, Lund-Hanssen H, Hofvind S. Possible strategies for use of artificial intelligence in screen-reading of mammograms, based on retrospective data from 122,969 screening examinations. *Eur Radiol.* 2022;32(12):8238-46.
15. Salim M, Wåhlin E, Dembrower K, Azavedo E, Foukakis T, Liu Y, et al. External Evaluation of 3 Commercial Artificial Intelligence Algorithms for Independent Assessment of Screening Mammograms. *JAMA Oncol.* 2020;6(10):1581-8.
16. Klompenhouwer EG, Duijm LE, Voogd AC, den Heeten GJ, Nederend J, Jansen FH, et al. Variations in screening outcome among pairs of screening radiologists at non-blinded double reading of screening mammograms: a population-based study. *Eur Radiol.* 2014;24(5):1097-104.
17. Brennan PC, Ganesan A, Eckstein MP, Ekpo EU, Tapia K, Mello-Thoms C, et al. Benefits of Independent Double Reading in Digital Mammography: A Theoretical Evaluation of All Possible Pairing Methodologies. *Acad Radiol.* 2019;26(6):717-23.
18. Gandomkar Z, Tay K, Brennan PC, Kozuch E, Mello-Thoms C. Can eye-tracking metrics be used to better pair radiologists in a mammogram reading task? *Med Phys.* 2018;45(11):4844-56.
19. Dembrower K, Lindholm P, Strand F. A Multi-million Mammography Image Dataset and Population-Based Screening Cohort for the Training and Evaluation of Deep Neural Networks-the Cohort of Screen-Aged Women (CSAW). *Journal of digital imaging.* 2020;33(2):408-13.
20. Taylor-Phillips S, Wallis MG, Parsons H, Dunn J, Stallard N, Campbell 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.
21. Taylor-Phillips S, Wallis MG, Jenkinson D, Adekanmbi V, Parsons H, Dunn J, et al. Effect of Using the Same vs Different Order for Second Readings of Screening Mammograms on Rates of Breast Cancer Detection: A Randomized Clinical Trial. *Jama.* 2016;315(18):1956-65.
22. Salim M, Dembrower K, Eklund M, Lindholm P, Strand F. Range of Radiologist Performance in a Population-based Screening Cohort of 1 Million Digital Mammography Examinations. *Radiology.* 2020;297(1):33-9.

23. Dembrower K, Wåhlin E, Liu Y, Salim M, Smith K, Lindholm P, et al. Effect of artificial intelligence-based triaging of breast cancer screening mammograms on cancer detection and radiologist workload: a retrospective simulation study. *Lancet Digit Health*. 2020;2(9):e468-e74.
24. Dembrower K, Liu Y, Azizpour H, Eklund M, Smith K, Lindholm P, et al. Comparison of a Deep Learning Risk Score and Standard Mammographic Density Score for Breast Cancer Risk Prediction. *Radiology*. 2020;294(2):265-72.
25. Taylor-Phillips S, Jenkinson D, Stinton C, Wallis MG, Dunn J, Clarke A. Double Reading in Breast Cancer Screening: Cohort Evaluation in the CO-OPS Trial. *Radiology*. 2018;287(3):749-57.
26. Cooper JA, Jenkinson D, Stinton C, Wallis MG, Hudson S, Taylor-Phillips S. Optimising breast cancer screening reading: blinding the second reader to the first reader's decisions. *Eur Radiol*. 2022;32(1):602-12.
27. Elin Wølner Bjørnson ÅSH, Silje Sagstad, Marthe Larsen, Jonas Thy, Gunhild Mangerud, Anne Kathrin, Ertzaas SH. *BreastScreen Norway: 25 years of organized screening*. Oslo: Cancer Registry of Norway; 2022.
28. Burnside ES, Lin Y, Munoz del Rio A, Pickhardt PJ, Wu Y, Strigel RM, et al. Addressing the challenge of assessing physician-level screening performance: mammography as an example. *PLoS One*. 2014;9(2):e89418.

Tables

Table 1 –Population characteristics and screening performance for the different datasets

	Sweden	England	Norway
Population characteristics			
Women screened, n	416,861	1,194,147	694,740
Median age at screening, years (IQR)	53 (46-62) ^A	59 (53-65) ^E	59 (54-64)
Screening performance			
Screening examinations, n	1,180,828	1,194,147	2,230,225
Recalls, per 1,000 examinations	23.0 ^B	41.5	34.2
CDR, per 1,000 examinations	3.4 ^{B, C}	8.8 ^F	5.9 ^F
Interval cancers, per 1,000 examinations	1.8 ^{A, B, D}	1.9 ^G	1.8 ^G

^A Age was unknown for 7 screening examinations

^B Final screening assessment was unknown for 106 screening examinations

^C Recall and breast cancer detected within 12 months after screening

^D Women who did not have a screen detected cancer, but had a breast cancer diagnosed within 18-24 months after screening

^E Age was unknown for 6 screening examinations

^F Breast cancer detected before the next screening examination as a result of recall at screening

^G Women who did not have a screen detected cancer, but had a breast cancer diagnosed before the next screening round. For the English dataset, interval cancers were incomplete, because not all data for interval cancers was known at the time of data extraction.

CDR, Cancer Detection Rate; IQR, interquartile range.

Table 2 – Population, reader, and pair characteristics for the study subsamples after selection criteria

	Sweden	England	Norway
	n=965,263 examinations	n=837,048 examinations	n=1,790,103 examinations
Population characteristics			
Women screened, n	396,193	837,048	647,275
Median age at screening, years (IQR)	53 (46-62)	59 (53-65) ^B	59 (54-64)
Reader characteristics			
Readers, n	36	326	121
Individual weighted mean AIR, per 1,000 examinations	36.3	48.2	49.0
Individual weighted mean CDR, per 1,000 examinations	3.1 ^A	8.3	5.2
Median cumulative reading volume, n (IQR)	27,346 (8,980-81,599)	4,524 (2,916-6,402)	20,259 (7,788-42,686)
Pair characteristics			
Pairs, n	64	560	203
Paired weighted mean AIR, per 1,000 examinations	47.8	63.7	77.0
Paired weighted mean CDR, per 1,000 examinations	3.4	8.9	6.0
Disagreement between two readers in a pair (%)	2.3	3.1	5.6
Median cumulative reading volume, n (IQR)	10,120 (5,529-17,119)	1,240 (938-1,727)	5,612 (3,445-10,344)

^A Recall and breast cancer diagnosed within 18-24 months after screening, 18 months for women ≤49 years old and 24 months for women >49 years old

^B Age was unknown for 5 screening examinations

AIR, Abnormal Interpretation Rate; CDR, Cancer Detection Rate; IQR, interquartile range.

Table 3 – Group screening performance for the different pairing strategies

	Sweden		England		Norway	
	AIR (95% CI)	CDR (95% CI)	AIR (95% CI)	CDR (95% CI)	AIR (95% CI)	CDR (95% CI)
Random	54.1 (46.1-62.1)	3.3 (2.8-3.9)	69.3 (63.7-74.9)	9.1 (8.3-9.9)	84.1 (80.3-88.0)	6.3 (5.9-6.7)
(Both) opposite	45.5 (19.2-71.7)	3.2 (2.4-4.1)	69.3 (65.7-72.9)	9.4 (6.7-12.2)	88.1 (74.8-101.3)	6.8 (5.3-8.3)
Opposite CDR	51.3 (37.5-65.0)	3.3 (2.4-4.3)	69.8 (66.5-73.1)	9.2 (7.7-10.6)	84.4 (77.1-91.6)	6.5 (5.7-7.4)
Opposite AIR	53.4 (38.1-68.7)	3.1 (2.6-3.6)	68.2 (61.2-75.2)	9.4 (7.9-10.8)	86.7 (80.6-92.7)	6.4 (5.7-7.1)
(Both) similar	52.3 (46.3-58.4)	3.6 (2.9-4.4)	70.5 (66.4-74.5)	8.9 (8.1-9.7)	82.4 (73.9-91.0)	6.1 (5.5-6.7)
Similar CDR	56.9 (46.6-67.1)	3.3 (2.9-3.8)	68.8 (65.2-72.3)	9.1 (8.3-9.9)	83.9 (78.0-89.7)	6.1 (5.5-6.6)
Similar AIR	54.7 (50.2-59.2)	3.6 (2.8-4.4)	70.3 (67.9-72.8)	8.9 (8.3-9.5)	81.6 (77.2-85.9)	6.2 (5.8-6.6)

AIR & CDR are given per 1,000 examinations and 95% CI are Bonferroni adjusted (P-values <0.05/6) confidence intervals, obtained by bootstrap resampling (n=1,000).

AIR, abnormal interpretation rate; CDR, cancer detection rate.

Figures

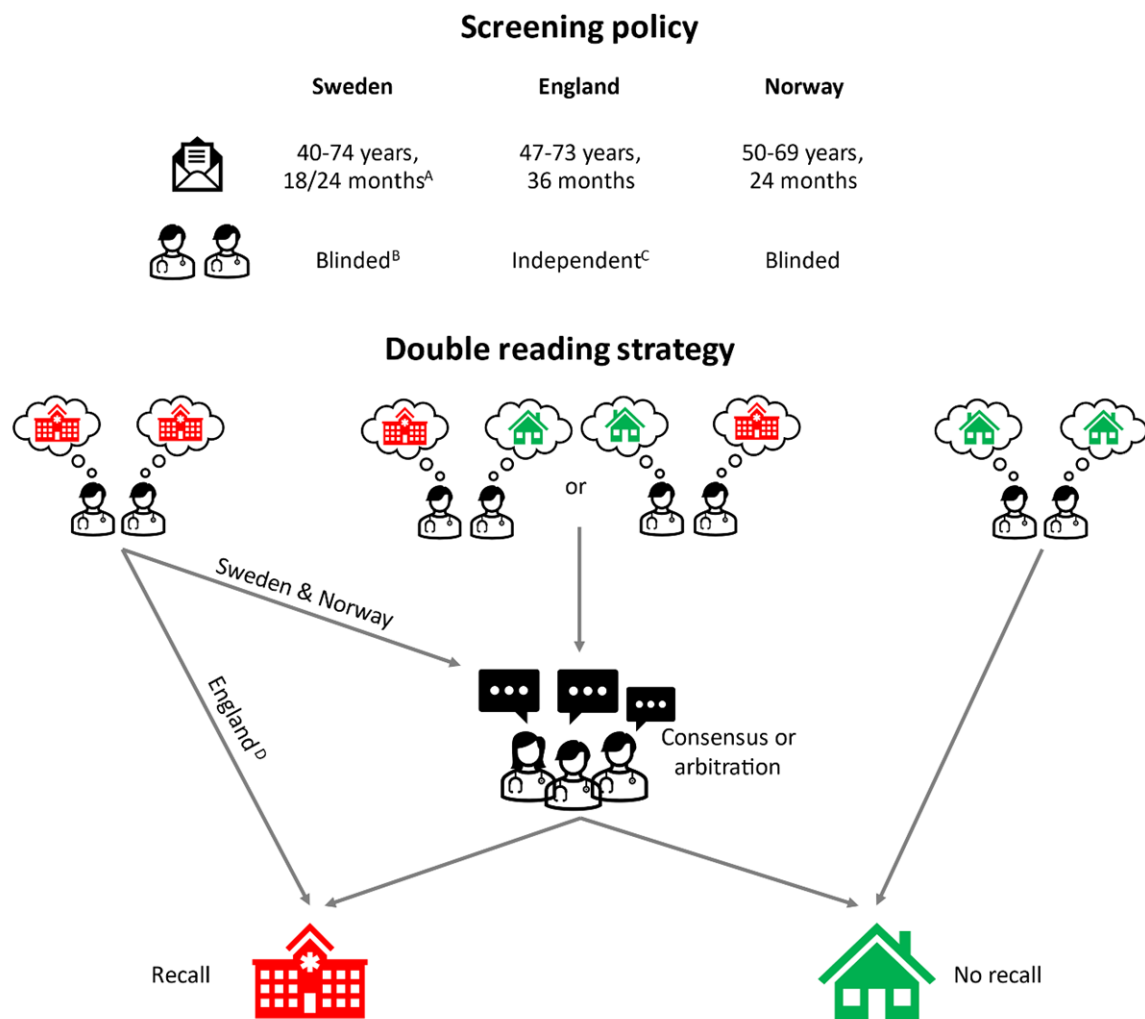


Figure 1 – Screening policy and double reading strategy for the different datasets.

^A Women >49 years old were invited every 24 months, whereas younger women were invited every 18 months; ^B There are local variations in practice regarding blinding, but at most centers the second reader was blinded to the decision of the first; ^C There are local variations in practice regarding blinding, but at most centers the second reader was not blinded to the decision of the first; ^D In most centers only discrepant readings were resolved through a single third reader or group arbitration, but in some centers with high recall rates, arbitration was also applied when both readers indicated recall.

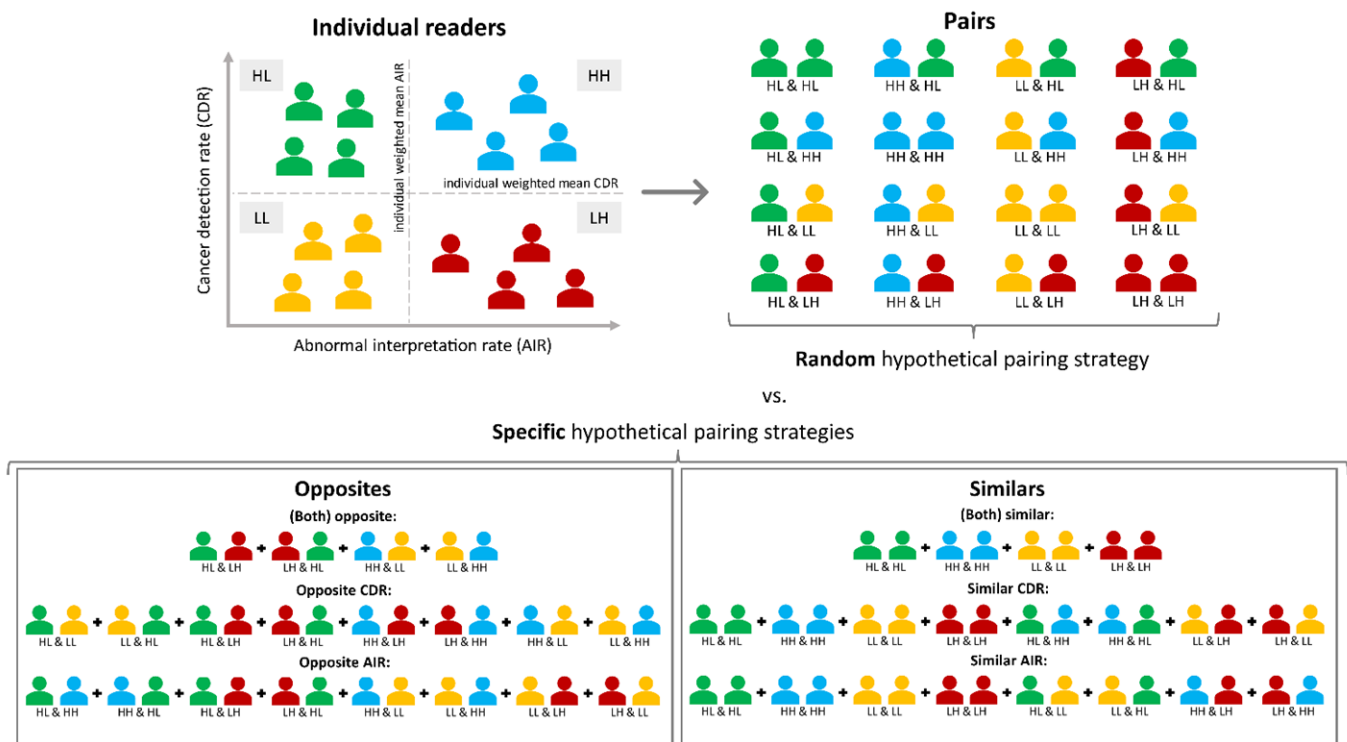


Figure 2 – Explanation for the evaluation of the screening performance of the individual readers and pairs, as well as an explanation of the hypothetical pairing strategies used to investigate the potential to improve outcomes with an appropriate double reading strategy.

AIR, abnormal interpretation rate; CDR, cancer detection rate; HH, high cancer detection rate & high abnormal interpretation rate; HL, high cancer detection rate & low abnormal interpretation rate; LH, low cancer detection rate & high abnormal interpretation rate; LL, low cancer detection rate & low abnormal interpretation rate.

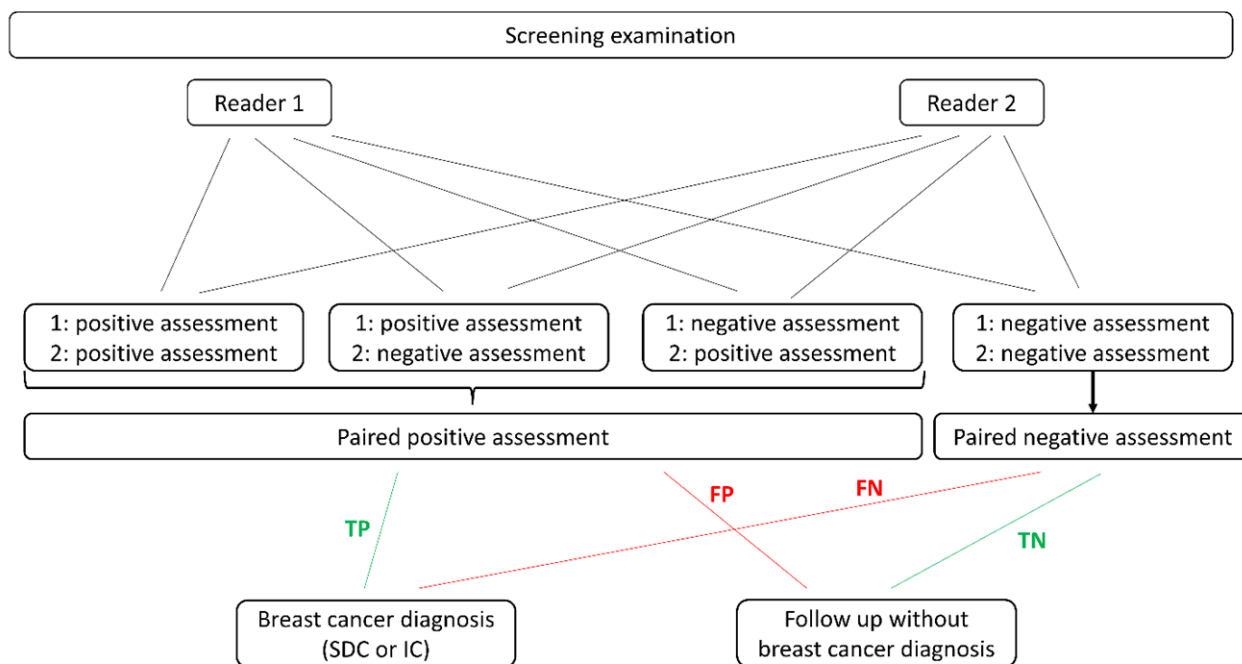


Figure 3 – Flowchart with the possible screening reading outcomes for the pairs in this study. Discrepant paired readings were defined as positive assessments. TP, True Positive; FP, False Positive; FN, False Negative; TN, True Negative; SDC, Screen-Detected Cancer; IC, Interval Cancer

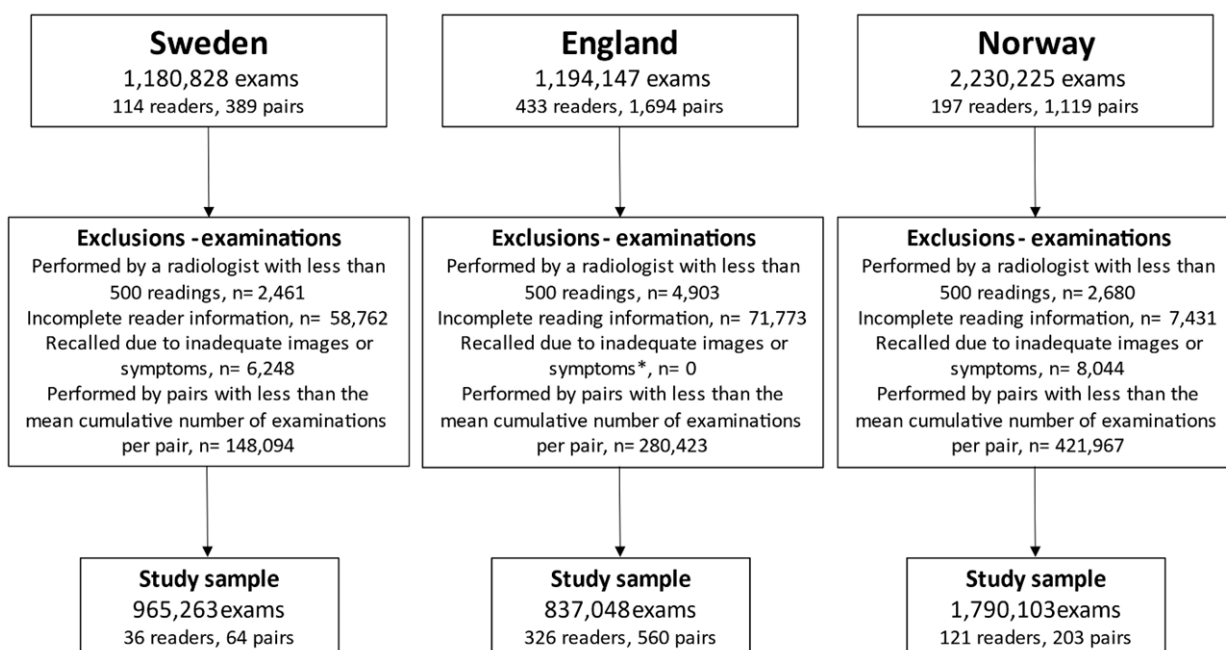


Figure 4 – Flowchart of screening examinations after applying exclusion criteria

* examinations of the English women who were recalled due to inadequate images or symptoms were already excluded before receiving the data

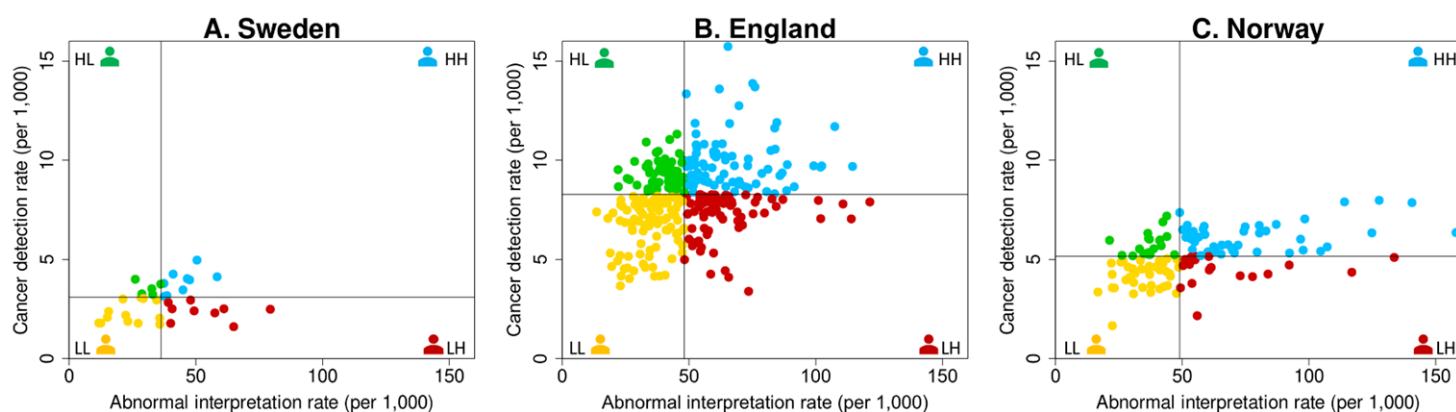


Figure 5 – Quadrant graphs for individual screening performance

Readers were classified into: high cancer detection & low abnormal interpretation (HL), high cancer detection & high abnormal interpretation (HH), low cancer detection & low abnormal interpretation (LL), or low cancer detection & high abnormal interpretation (LH). The individual weighted mean performance was used as cutoff.

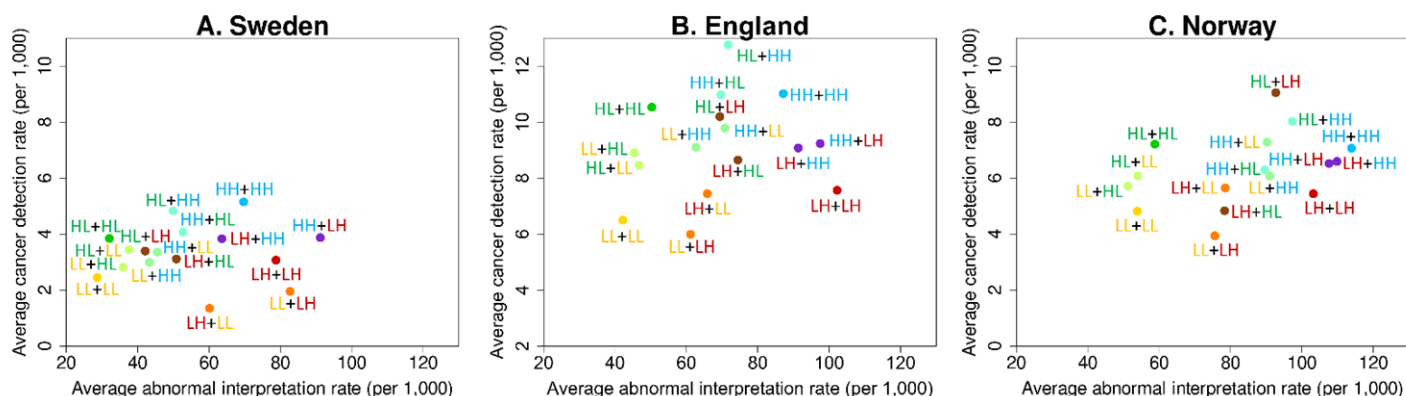


Figure 6 – Scatter plots with the average CDR and AIR of the sixteen specific pairs.

HH, high cancer detection rate & high abnormal interpretation rate; HL, high cancer detection rate & low abnormal interpretation rate; LH, low cancer detection rate & high abnormal interpretation rate; LL, low cancer detection rate & low abnormal interpretation rate.

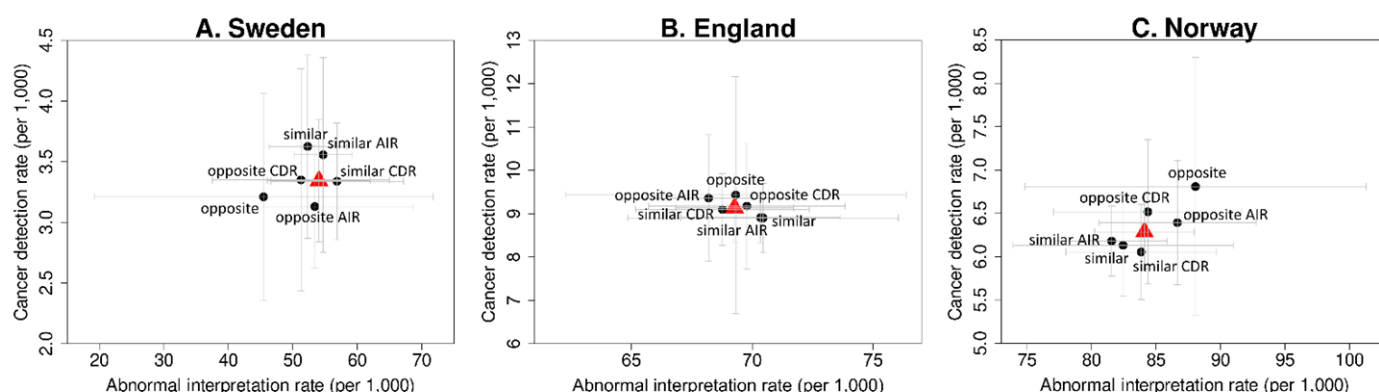


Figure 7 – Group screening performance for the different pairing strategies.

The triangles (red) represent the average screening performance for the random hypothetical pairing strategy and the dots represent the performance for the specific hypothetical pairing strategies (black). Error bars are Bonferroni adjusted 95% confidence intervals, obtained by bootstrap resampling ($n=1,000$). Please note that the axes are different, due to the differences in CDR and AIR for the datasets. AIR, abnormal interpretation rate; CDR, cancer detection rate.