

1 **Reliability of single-lead electrocardiogram interpretation to detect atrial**
2 **fibrillation: insights from the SAFER Feasibility Study**

3 Katie Hibbitt,^a James Brimicombe,^a Martin R. Cowie,^b Andrew Dymond,^a Ben Freedman,^c
4 Simon J Griffin,^a FD Richard Hobbs,^f Hannah Clair Lindén,^c Gregory Y. H. Lip,^d Jonathan
5 Mant,^a Richard J. McManus,^f Madhumitha Pandiaraja,^a Kate Williams,^a Peter H. Charlton,^{a*}
6 on behalf of the SAFER Investigators

7 ^aDepartment of Public Health and Primary Care, University of Cambridge, Cambridge, CB1
8 8RN, UK

9 ^bRoyal Brompton Hospital (Guy's and St Thomas' NHS Foundation Trust), Sydney Street,
10 London, SW3 6NP, UK

11 ^cZenicor Medical Systems AB, 113 59 Stockholm, Sweden

12 ^dLiverpool Centre for Cardiovascular Science at University of Liverpool, Liverpool
13 John Moores University and Liverpool Heart & Chest Hospital, Liverpool, United
14 Kingdom; and Danish Center Health Services Research, Department of Clinical
15 Medicine, Aalborg University, Aalborg, Denmark

16 ^eHeart Research Institute, University of Sydney, Sydney 2006, Australia

17 ^fNuffield Department of Primary Care Health Sciences, University of Oxford, Oxford, OX2
18 6GG, UK

19 * PHC corresponding author.

20

21 **Corresponding author:** Dr Peter H Charlton

22 Department of Public Health and Primary Care, University of Cambridge, Cambridge

23 Telephone number: 01223 331063

24 Email: pc657@cam.ac.uk

25

NOTE: This preprint reports new research that has not been certified by peer review and should not be used to guide clinical practice.

26 **Funding:** This study is funded by the National Institute for Health Care Research (NIHR),
27 grant number RP-PG-0217-20007; and the British Heart Foundation (BHF), grant number
28 FS/20/20/34626. The views expressed are those of the authors and not necessarily those of the
29 NIHR or the Department of Health and Social Care.

30

31 **Disclosures:** MRC is employed by Astrazeneca PLC. BF has received speaker fees, honoraria,
32 and non-financial support from the BMS and Pfizer Alliance; and loan devices for investigator
33 initiated studies from Alivecor: all were unrelated to the present study, but related to screening
34 for AF. SJG has received honoraria from Astra Zeneca and Eli Lilly for contributing to
35 postgraduate education concerning type 2 diabetes to specialists and primary care teams. FDRH
36 reports occasional consultancy for BMS/Pfizer, Bayer and BI over the last five years. HCL is
37 employed by Zenicor Medical Systems AB. GYHL: Consultant and speaker for BMS/Pfizer,
38 Boehringer Ingelheim, Daiichi-Sankyo, Anthos. No fees are received personally. JM has
39 performed consultancy work for BMS/Pfizer and Omron. PHC has performed consultancy
40 work for Cambridge University Technical Services, and has received honoraria from IOP
41 Publishing and Emory University (the latter not received personally).

42

43 **Data availability:** Requests for pseudonymised data should be directed to the SAFER study
44 co-ordinator (Andrew Dymond using SAFER@medschl.cam.ac.uk) and will be considered by
45 the investigators, in accordance with participant consent.

46

47 **Total word count:** 3,838 (exc. acknowledgments, references, figure legends and tables)

48 **Number of tables:** 3

49 **Number of figures:** 3

50

51 **Abstract**

52 **Background**

53 Single-lead electrocardiograms (ECGs) can be recorded using widely available devices such
54 as smartwatches and handheld ECG recorders. Such devices have been approved for atrial
55 fibrillation (AF) detection. However, little evidence exists on the reliability of single-lead ECG
56 interpretation. We aimed to assess the level of agreement on detection of AF by independent
57 cardiologists interpreting single lead ECGs, and to identify factors influencing agreement.

58 **Methods**

59 In a population-based AF screening study, adults aged ≥ 65 years old recorded four single-lead
60 ECGs per day for 1-4 weeks using a handheld ECG recorder. ECGs showing signs of possible
61 AF were identified by a nurse with the aid of an automated algorithm. These ECGs were
62 reviewed by two independent cardiologists who assigned participant- and ECG-level
63 diagnoses. Inter-rater reliability of AF diagnosis was calculated using linear weighted Cohen's
64 kappa (κ_w).

65 **Results**

66 185 participants and 1,843 ECGs were reviewed by both cardiologists. The level of agreement
67 was moderate: $\kappa_w = 0.42$ (95% CI, 0.32 – 0.52) at the participant-level; and $\kappa_w = 0.51$ (0.46 –
68 0.56) at the ECG-level. At participant-level, agreement was associated with the number of
69 adequate-quality ECGs recorded, with higher agreement in participants who recorded at least
70 67 adequate-quality ECGs. At ECG-level, agreement was associated with ECG quality and
71 whether ECGs exhibited algorithm-identified possible AF.

72 **Conclusions**

73 Inter-rater reliability of AF diagnosis from single-lead ECGs was found to be moderate in older
74 adults. Strategies to improve reliability might include participant and cardiologist training and
75 designing AF detection programmes to obtain sufficient ECGs for reliable diagnoses.

76

77 **Clinical Trial Registration:** ISRCTN 16939438; <https://doi.org/10.1186/ISRCTN16939438>

78

79 **Keywords**

80 Atrial fibrillation, diagnosis, electrocardiogram, inter-rater agreement, screening

81

82

83

84

85

86

87

88

89 1. Introduction

90 The electrocardiogram (ECG) is a fundamental technique for assessing the functionality
91 of the heart. The process for recording a 12-lead ECG was described 70 years ago (1), and to
92 this day the 12-lead ECG remains widely used for the diagnosis and management of a range of
93 heart conditions (2). Whilst the 12-lead ECG is highly informative, providing several ‘views’
94 of the heart’s electrical activity, it can only be measured by clinicians in a healthcare setting.
95 Recently, clinical and consumer devices have become available which allow individuals to
96 record a single-lead ECG on demand via a smartwatch or handheld device. This approach has
97 a number of useful features: such ECGs can be measured by patients themselves with no
98 clinical input, can be acquired synchronously with symptoms, can be repeated on multiple
99 occasions with minimal inconvenience and can be transmitted electronically to healthcare
100 providers (3).

101 Atrial fibrillation (AF) is a common arrhythmia which confers a fivefold increase in the
102 risk of stroke (4) which can be mitigated through anticoagulation (5). A significant proportion
103 of AF remains unrecognised (6) as it may be asymptomatic or occur only intermittently. Self-
104 captured, single-lead ECGs could greatly assist in the detection of AF (7) when: (i) used by
105 device owners, with ECGs acquired opportunistically, upon symptoms, or when prompted by
106 a device (8); and when (ii) used in screening programmes, allowing multiple ECGs to be
107 acquired from an individual over a period of weeks (9). Indeed, European Society of
108 Cardiology guidelines support the use of single-lead ECGs acquired from wearable or mobile
109 devices to identify AF (8). Whilst automated algorithms can be used to identify those ECGs
110 which show evidence of AF and therefore warrant clinical review (10), a final diagnosis of AF
111 must be made by a physician interpreting an ECG (8). To date, there is little evidence on the
112 reliability of single-lead ECG interpretation for AF diagnosis, and most existing evidence is
113 derived from ECGs collected from hospital patients (11-13).

114 We aimed to assess the level of agreement on detection of AF by independent
115 cardiologists interpreting single lead ECGs, and to identify factors which influence agreement.

116 2. Methods

117 We assessed inter-rater agreement using ECG data collected in a population-based AF
118 screening study, in which each participant recorded multiple ECGs. Agreement between
119 cardiologist interpretations was assessed at the participant-level (*i.e.* the overall participant
120 diagnosis) and the ECG-level (*i.e.* interpretations of individual ECGs). In addition, we
121 investigated the influence of several factors on the level of agreement (*e.g.* participant age and
122 ECG quality).

123 2.1 Data collection

124 We collected the data for these analyses in the SAFER (Screening for Atrial Fibrillation
125 with ECG to Reduce stroke) Feasibility Study (ISRCTN 16939438), conducted in 2019 and
126 approved by the London-Central Research Ethics Committee (REC ref: 18/LO/2066).
127 Participants were older adults aged ≥ 65 years old, who were not receiving long-term
128 anticoagulation for stroke prevention, not on the palliative care register, and not resident in a
129 nursing home. All participants gave written informed consent.

130 In this study, older adults (aged 65 and over) recorded single-lead ECGs at home using
131 the handheld Zenicor EKG-2 device (Zenicor Medical Systems AB) (10). This device measures
132 a 30-second, single-lead ECG between the thumbs, using dry electrodes. Participants were
133 asked to record four ECGs per day for either 1, 2 or 4 weeks. The ECGs were transferred to a
134 central database for analysis and review.

135 Participant- and ECG-level diagnoses were obtained as follows (and as summarised in
136 Figure 1). First, a computer algorithm was used to identify abnormal ECGs (Cardiolund ECG
137 Parser algorithm, Cardiolund AB, Sweden). The algorithm has previously been found to have
138 a sensitivity for AF detection of approximately 98% (10). Second, a nurse reviewed all the
139 ECGs which were classified by the algorithm as abnormal, and manually corrected any

140 algorithm misclassifications based on their clinical judgement. The nurse then identified
141 participants for cardiologist review as those participants with at least one ECG classified as
142 abnormal which the nurse deemed exhibited signs of possible AF (as detailed in (14)). Third,
143 these participants were sent for review by two highly experienced cardiologists, both of whom
144 had substantial ECG reviewing experience (GYHL and MRC). The cardiologists had access to
145 all the ECGs from these participants, though it was not anticipated that ECGs that were
146 classified as normal would be reviewed, or that all abnormal ECGs would be reviewed, once a
147 participant-level diagnosis had been reached. Each cardiologist independently provided a
148 diagnosis for each participant. For AF to be diagnosed it was required to be present for the
149 whole 30 seconds, or the entire trace where the ECG was interpretable. No other formal
150 definition of AF was provided for the cardiologists to use. In addition, on an *ad hoc* basis, the
151 cardiologists also provided diagnoses for individual ECGs and labelled ECGs as ‘low-quality’.
152 Diagnoses were categorised as: AF \geq 30 seconds duration; cannot exclude AF; or, non-AF.

153 We extracted a subset of the collected data for the analysis as follows. Only data from
154 those participants who were reviewed by both cardiologists were included in participant-level
155 analyses. In addition, only those ECGs which were reviewed by both cardiologists were
156 included in ECG-level analyses. We excluded from analyses any ECGs for which a
157 cardiologist’s initial diagnosis was not recorded (prior to subsequent resolution of
158 disagreements).

159 **2.2 Data Processing**

160 We obtained the characteristics of each ECG as follows: First, the computer algorithm
161 extracted the following characteristics: heart rate, ECG quality (either normal or poor quality),
162 level of RR-interval variability (calculated as the standard deviation of RR-intervals divided
163 by the mean RR-interval), and whether or not an ECG exhibited algorithm-identified possible

164 AF (defined as the ECG having either irregular RR-intervals or a fast regular heart rate).
165 Second, the quality of ECGs was obtained by combining the quality assessment provided by
166 the algorithm with cardiologist comments on ECG quality: any ECGs which the algorithm or
167 at least one cardiologist deemed to be of poor quality were classed as low quality in the analysis.
168 ECGs for which the algorithm was unable to calculate heart rate or RR-interval variability were
169 excluded from analyses requiring those characteristics.

170 **2.3 Statistical Analysis**

171 We assessed the reliability of ECG interpretation using both participant-level diagnoses
172 and ECG-level cardiologist diagnoses. First, we reported the overall levels of agreement.
173 Second, we assessed the influence of different factors on levels of agreement, such as the
174 influence of ECG quality. The factors assessed at the participant-level were: age, gender,
175 number of adequate-quality ECGs recorded by a participant, and the number of ECGs recorded
176 by a participant exhibiting algorithm-identified possible AF. The factors assessed at the ECG-
177 level were: heart rate, RR-interval variability, ECG quality, and whether or not an ECG
178 exhibited algorithm-identified possible AF. We investigated factors which were continuous
179 variables (such as heart rate) by grouping values into categories with similar sample sizes (*e.g.*
180 heart rates were categorised as 30-59 bpm, 60-69 bpm, etc).

181 We assessed agreement between cardiologists using inter-rater reliability statistics. The
182 primary statistic, Cohen's kappa, κ , provides a measure of the difference between the actual
183 level of agreement between cardiologists, and the level of agreement that would be expected
184 by random chance alone. Values for κ range from -1 to 1, with -1 indicating complete
185 disagreement; 0 the level expected by chance; 0.01-0.20 slight agreement; 0.21-0.40 fair
186 agreement; 0.41-0.60 moderate agreement; 0.61-0.80 substantial agreement; 0.81-0.99 almost
187 perfect agreement; and 1 perfect agreement (15). The second statistic, a weighted Cohen's
188 kappa, κ_w , reflects the greater consequences of a disagreement of 'AF' vs 'non-AF', compared

189 to a disagreement of ‘cannot exclude AF’ vs either ‘AF’ or ‘non-AF’. We weighted
190 disagreements of ‘AF’ vs ‘non-AF’ as complete disagreements, whereas disagreements
191 including ‘cannot exclude AF’ were weighted equivalently to the level expected by chance.
192 We reported the third statistic, percentage agreement, to facilitate comparisons with previous
193 studies.

194 We calculated 95% confidence intervals for κ and κ_w using bootstrapping. We
195 undertook tests for significant associations between factors (*e.g.* heart rate) and the level of
196 agreement using a chi-square test for independence between the proportion of agreement in
197 each category.

198 **3. Results**

199 Of the 2,141 participants who were screened, 190 had ECGs which underwent
200 cardiologist review and were therefore included in the participant-level analyses, as shown in
201 Figure 1. Most participants' ECGs were not sent for cardiologist review (1,951 participants)
202 because either: (i) the computer algorithm did not find any abnormalities in their ECGs (603
203 participants); or (ii) the nurse reviewer judged that none of their abnormal ECGs exhibited
204 signs of possible AF (1,348 participants).

205 [Figure 1]

206 The 190 participants whose ECGs underwent cardiologist review recorded a total of
207 15,258 ECGs, with a median (lower – upper quartiles) of 67.0 (56.0 - 112.0) ECGs each. The
208 two cardiologists assigned diagnoses to 1,996 and 4,411 of these ECGs respectively of which
209 1,872 ECGs were assigned diagnoses by both cardiologists. Initial diagnoses (prior to
210 subsequent resolution of disagreements) were not recorded for 29 of these ECGs, leaving 1,843
211 available for ECG-level analyses (see Figure 2).

212 [Figure 2]

213 **3.1. Reliability of AF diagnosis at the participant-level**

214 The inter-rater reliability of AF diagnosis at the participant-level, when the
215 cardiologists had access to all the ECGs recorded by a participant, was moderate ($\kappa_w =$
216 $0.48 (0.37 - 0.58)$; $\kappa = 0.42 (0.32 - 0.52)$; and $\%_{agree} = 66.3\%$ (

217 Table 1).

218 [

219 Table 1]

220 The results for the relationship between the level of agreement between cardiologists
221 and factors at the participant- and ECG-level are presented in Figure 3 and Table 2. At the
222 participant-level, the level of agreement was significantly associated with the number of
223 adequate-quality ECGs recorded by a participant. Participants who recorded at least 67
224 adequate-quality ECGs had a significantly higher level of agreement in their diagnoses than
225 those who recorded fewer than 67. There was agreement on 52.6% of participant-level
226 diagnoses in those participants with <67 adequate-quality ECGs, compared to 80.0% in those
227 with 67 or more. Of the 31 participants for whom there was complete disagreement (where one
228 cardiologist diagnosed AF and the other diagnosed non-AF), 23 (74%) recorded <67 adequate-
229 quality ECGs. There was no significant association between the level of agreement and age or
230 gender.

231 [Table 2]

232 [Figure 3]

233 3.2. Reliability of ECG interpretation

234 The inter-rater reliability of AF diagnosis at the individual ECG-level was moderate
235 ($\kappa_w = 0.58$ (0.53 – 0.63); $\kappa = 0.51$ (0.46 – 0.56); and $\%_{agree} = 86.1\%$ (Table 3).
236 Referring to the ECG-level results in Figure 3 and Table 2, the level of agreement was
237 significantly associated with ECG quality, with low-quality ECGs associated with a lower level
238 of agreement. This remained regardless of whether quality was assessed using cardiologist
239 comments on ECG quality, the automated algorithm assessment, or a combination of both. The
240 level of agreement was also significantly associated with whether or not an ECG exhibited
241 algorithm-identified possible AF, where ECGs exhibiting possible AF were associated with a

242 higher level of agreement. There was no significant association between the level of agreement
243 and heart rate or RR-interval variability.

244 [Table 3]

245 3.3. Comparison of cardiologists' reviewing practices

246 The two cardiologists' reviewing practices differed. At the participant-level, one
247 cardiologist diagnosed more participants with AF than the other (72 out of 190, *i.e.* 38%, vs.
248 50, *i.e.* 26%) (see Table 1). Similarly, at the ECG-level, this cardiologist diagnosed more ECGs
249 as AF than the other (235 out of 1,843, *i.e.* 13%, vs. 179, *i.e.* 9.7%), and more ECGs as 'cannot
250 exclude AF' than the other (119, *i.e.* 6%, vs. 63, *i.e.* 3%) (see Table 3). Most of the ECGs
251 diagnosed as AF by the cardiologists exhibited an irregular rhythm as identified by the
252 algorithm (95% of the 235 ECGs diagnosed as AF by one cardiologist, 88% of the 179 ECGs
253 diagnosed as AF by the other cardiologist, and 95% of the 137 ECGs diagnosed as AF by both
254 cardiologists).

255

256 **4. Discussion**

257 **4.1. Summary of findings**

258 This study provides evidence on the inter-rater reliability of single-lead
259 electrocardiogram interpretation, and the factors that influence this. Moderate agreement was
260 observed between cardiologists on participant-level diagnoses of AF in a population-based AF
261 screening study when this diagnosis was made using multiple ECGs per participant. The key
262 factor associated with the level of agreement at the participant-level was the number of
263 adequate-quality ECGs recorded by a participant, with higher levels of agreement in those who
264 recorded more adequate-quality ECGs. Moderate agreement was observed between
265 cardiologists on the diagnoses of individual ECGs. Similarly, at the ECG-level, low-quality
266 ECGs were associated with lower levels of agreement. In addition, lower levels of agreement
267 were observed on those ECGs not exhibiting algorithm-identified possible AF.

268 **4.2. Comparison with existing literature**

269 The levels of agreement in AF diagnosis from single-lead ECGs observed in this study
270 are lower than in many previous studies. Previous studies have found almost perfect agreement
271 when interpreting 12-lead ECGs, but lower levels of agreement when interpreting single-lead
272 ECGs. In an analysis of 12-lead ECGs from the SAFE AF Screening Trial, cardiologists agreed
273 on the diagnosis of 99.7% of ECGs (all but 7 of 2,592 analysed ECGs)(16). In comparison, in
274 the present study of single-lead ECGs cardiologists agreed on the diagnosis of 86.1% of ECGs
275 (1,587 out of 1,843 ECGs). However, the proportion of normal ECGs included in this study
276 was substantially lower than in the SAFE AF Screening Trial (less than 1% in this study, versus
277 93% in SAFE), so the simple level of agreement is not directly comparable. Similarly, in a
278 study of the diagnosis of supraventricular tachycardia in hospital patients, an almost perfect
279 agreement of $\kappa = 0.97$ was observed in interpretation of 12-lead ECGs, compared to a

280 substantial agreement of $\kappa = 0.76$ when using single-lead ECGs from the same patients (17).
281 The previously reported levels of agreement for the diagnosis of AF from single-lead ECGs
282 have varied greatly between studies: fair agreement was observed by Kearley *et al.* (18) ($\kappa =$
283 0.28); moderate agreement was observed by Lowres *et al.* (19) (weighted $\kappa = 0.4$); substantial
284 agreements were observed by Poulsen *et al.* (12) ($\kappa = 0.65$) and Kearley *et al.* (18) ($\kappa = 0.76$);
285 and almost perfect agreements were observed by Desteghe *et al.* (11) ($\kappa = 0.69$ to 0.86), Koshy
286 *et al.* (20) ($\kappa = 0.80$ to 0.83), Wegner *et al.* (13) ($\kappa = 0.90$), and Racine *et al.* (21) ($\kappa = 0.94$).
287 The variation in levels of agreement may have been contributed to by study setting and
288 underlying frequency of AF, since those studies which reported the lowest levels of agreement
289 took place out-of-hospital (18, 19). The present study, conducted in the community, similarly
290 observed lower levels of agreement than many other studies ($\kappa = 0.42$ at the participant-level,
291 and $\kappa = 0.51$ at the ECG-level). In the context of AF screening, a 69.2% level of agreement
292 has been reported in a previous AF screening study by Pipilas *et al.* (22), compared to 86.1%
293 in the present screening study. The low levels of agreement in the present study could have
294 been contributed to by: (i) the ECGs being more challenging to review as an algorithm and a
295 nurse filtered out most ECGs which did not exhibit signs of AF (and are therefore easier to
296 interpret) prior to cardiologist review; (ii) the ECGs being of lower quality since participants
297 recorded ECGs themselves without clinical supervision; and (iii) the use of an additional
298 diagnostic category of ‘cannot exclude AF’.

299 This study’s findings about factors which influence the reliability of ECG interpretation
300 complement those reported previously (11,12,13,22,23). It has previously been reported that
301 ECGs exhibiting baseline wander, noise, premature beats, or low-amplitude atrial activity are
302 associated with mis-diagnoses (11,23). In this study low-quality ECGs were similarly
303 associated with lower levels of agreement between cardiologists. The significant proportion of
304 low-quality ECGs obtained when using a handheld ECG device has been reported previously,

305 with 12% of ECGs being judged as ‘very low quality’ in (22), 13% as ‘not useable’ in (12),
306 and 20% as ‘inadequate quality’ in (13).

307 The accuracy of both automated and manual diagnosis of AF from single-lead ECGs
308 has been assessed previously. A recent meta-analysis found pooled sensitivities and
309 specificities of automated ECG diagnoses of 89% and 99% respectively in the community
310 setting (24). The accuracy of manual diagnoses has varied greatly between previous studies,
311 with sensitivities and specificities in comparison to reference 12-lead ECGs reported as: 77.4%
312 and 73.0% (22), 90% and 79% (21), 76-92% and 84-100% (20), 89-100% and 85-88% (25),
313 92.5% and 89.8% (26), 93.9% and 90.1% (18), 100% and 94% (13), and 100% and 100% (27).
314 In all of these studies, the single-lead ECGs were recorded under supervision. In contrast, the
315 present study considered ECGs collected using a telehealth device at home without supervision.

316 **4.3. Strengths and limitations**

317 There are several strengths to this study. First, we assessed the level of agreement in
318 both participant-level ECG-level AF diagnoses which is of particular relevance in AF
319 screening, whereas most previous work has been limited to ECG-level diagnoses. Second, the
320 ECGs used in this study were collected in a prospective population-based AF screening study,
321 and are therefore representative of ECGs captured in telehealth settings by older adults without
322 clinical supervision. The ECGs were recorded using dry electrodes, as opposed to the gel
323 electrodes used in clinical settings. Dry electrodes can result in poorer conduction and therefore
324 lower signal quality, making interpretation more challenging. Since smartwatches also use dry
325 electrodes, the findings are expected to be relevant to the growing use of ECG-enabled
326 consumer devices. Third, the ECGs included in the analysis are representative of those which
327 would be sent for clinical review in real-world settings: ECGs without signs of abnormalities
328 were excluded using an automated, CE-marked analysis system, leaving only those ECGs with

329 signs of abnormalities for review. Fourth, the study included a large number of ECGs (1,843),
330 each interpreted by two cardiologists. Fifth, we used Cohen's kappa statistic to assess the level
331 of agreement between cardiologists: this statistic takes into account agreement by chance
332 unlike the percentage agreement (Viera and Garrett, 2005).

333 The key limitations to this study are as follows. First, the findings are based on data
334 from only 190 participants who had an abnormal ECG flagged in the study. Second, inter-rater
335 agreement was assessed using diagnoses provided by only two cardiologists. Third, not all
336 ECGs sent for review were interpreted by both cardiologists, with those not interpreted by both
337 cardiologists excluded from the analysis. Fourth, the initial diagnosis was not recorded for a
338 small minority of the ECGs reviewed by both cardiologists (29 out of 1,872, 1.5%), so these
339 were not included in the analysis. Finally, it should be remembered that the study assessed the
340 reliability of ECG interpretation (*i.e.* the level of agreement between two cardiologists), rather
341 than the accuracy of ECG interpretation (*i.e.* a comparison of cardiologist interpretation against
342 an independent reference). In doing so, the study identified factors associated with reduced
343 levels of agreement, providing evidence on how to improve the level of agreement, and
344 subsequently the reliability of interpretation.

345 **4.4. Implications**

346 This study indicates that steps should be taken to ensure diagnoses based on single-lead
347 ECGs are as reliable as possible. Out of 2,141 participants screened for AF, there was
348 agreement between cardiologists on diagnoses of AF for 44 participants, complete
349 disagreement for 31 participants (AF vs. non-AF), and partial disagreement for 33 participants
350 (AF or non-AF vs. cannot exclude AF). In terms of disease prevalence, there was agreement
351 on AF diagnosis in 2.1% of the sample population, complete disagreement in 1.4%, and partial
352 disagreement in 1.5%.

353 The findings could inform the design of AF screening programmes. AF screening
354 programmes often include collection of multiple short ECGs (or a continuous ECG recording)
355 over a prolonged period to capture even infrequent episodes of paroxysmal AF. The results of
356 this study indicate that a prolonged period is also required to obtain reliable diagnoses: at least
357 67 adequate-quality ECGs were required for a reliable diagnosis in this study, providing
358 evidence that screening programmes should be designed to capture at least this many adequate-
359 quality ECGs from all participants (*i.e.* at least 17 days of screening when recording 4 ECGs
360 per day, and potentially 21 days of screening to account for missed or low-quality ECGs). In
361 addition, no association was found between participant gender or age and the reliability of
362 diagnoses, indicating that it is reasonable to use single-lead ECGs in older adults of a wide
363 range of ages (from 65 to 90+ in this study).

364 The findings of this study could also underpin strategies to obtain more reliable
365 participant-level diagnoses through personalised screening. Those individuals who are likely
366 to receive a less reliable diagnosis could be identified by using an automated algorithm to
367 analyse the quality of incoming ECGs, and then the duration of screening could be extended in
368 those individuals without sufficient adequate-quality ECGs. This could help increase reliability
369 by increasing the number of adequate-quality ECGs available for diagnosis. Second,
370 participants with a high proportion of low-quality ECGs could be offered additional training
371 on ECG measurement technique, potentially by telephone.

372 This study highlights the need to ensure single-lead ECG interpreters receive sufficient
373 training. The ECGs in this study were interpreted by highly experienced cardiologists, and yet
374 there was still disagreement over diagnoses for 16% of those participants sent for cardiologist
375 review. If single-lead ECG-based AF screening is widely adopted in the future, then it will be
376 important to ensure all ECG interpreters receive sufficient training and gain sufficient
377 experience in single-lead ECG interpretation to provide reliable diagnoses. We note that single-

378 lead ECG interpretation presents additional challenges beyond those encountered in 12-lead
379 ECG interpretation: ECGs may be of lower quality (12), P-waves may not be as visible (11),
380 and only one lead is available.

381 The findings of this study indicate that it is important that the quality of single-lead
382 ECGs is as high as possible. Particularly given the implications of an AF diagnosis such as
383 recommendations for anticoagulation treatment which increases the risk of bleeding. The
384 development of consumer and telehealth ECG devices involves making a range of design
385 decisions which can influence the quality of ECGs sent for clinical review, including: the size,
386 type, and anatomical position of electrodes; the filtering applied to signals to reduce noise; and
387 whether to exclude ECGs of insufficient quality from clinical review (and if so, how best to
388 identify these ECGs). Device should designers consider the potential effects of these design
389 decisions on the reliability of diagnoses.

390 **5. Conclusion**

391 Moderate agreement was found between cardiologists when diagnosing AF from
392 single-lead ECGs in an AF screening study. The study indicates that for every 100 screening
393 participants diagnosed with AF by two cardiologists, there would be complete disagreement
394 over the diagnosis of 70 further participants. This provides great incentive for ensuring that the
395 interpretation of single-lead ECGs is as reliable as possible. Key factors were identified which
396 influence the reliability of single-lead ECG interpretation. Most importantly, the quality of
397 ECG signals greatly influenced reliability. In addition, when multiple ECGs were acquired
398 from an individual, the reliability of participant-level diagnoses was influenced by the number
399 of adequate-quality ECGs available for interpretation. This new evidence could help improve
400 single-lead ECG interpretation, and consequently increase the effectiveness of screening for
401 AF using single-lead ECG devices. Future work should investigate how to obtain ECGs of the
402 highest possible quality in the telehealth setting, and how best to train ECG interpreters to
403 ensure diagnoses are as accurate as possible.

404

405 **Acknowledgment**

406 This study is funded by the National Institute for Health and Care Research (NIHR),
407 Programme Grants for Applied Research Programme (Reference Number RP-PG0217-20007);
408 the NIHR School for Primary Care Research (SPCR-2014-10043, project 410); and the British
409 Heart Foundation (BHF) grant number FS/20/20/34626. The views expressed are those of the
410 authors and not necessarily those of the NIHR or the Department of Health and Social Care.

411

412 References

- 413 1. Wilson FN, Kossmann CE, Burch GE, et al. Recommendations for Standardization of
414 Electrocardiographic and Vectorcardiographic Leads. *Circulation*. 1954;10(4):564-573.
- 415 2. Rafie N, Kashou AH, Noseworthy PA. ECG Interpretation: Clinical Relevance,
416 Challenges, and Advances. *Hearts*. 2021;2(4):505-513.
- 417 3. Hall A, Mitchell ARJ, Wood L, Holland C. Effectiveness of a single lead AliveCor
418 electrocardiogram application for the screening of atrial fibrillation. *Medicine (Baltimore)*.
419 2020;99(30):e21388.
- 420 4. Wolf PA, Abbot RD, Kannel WB. Atrial fibrillation as an independent risk factor for
421 stroke: the Framingham study. *Stroke*. 1991;22:983–8.
- 422 5. Ruff CT, Giugliano RP, Braunwald E, et al. Comparison of the efficacy and safety of
423 new oral anticoagulants with warfarin in patients with atrial fibrillation: a meta-analysis of
424 randomised trials. *The Lancet*. 2014;383(9921):955-962.
- 425 6. Public Health England. Atrial fibrillation prevalence estimates in England: Application
426 of recent population estimates of AF in Sweden. PHE Publications Gateway Number:
427 2014778. PHE Publications Gateway Number 2014778; 2015.
- 428 7. Svennberg E, Engdahl J, Al-Khalili F, Friberg L, Frykman V, Rosenqvist M. Mass
429 screening for untreated atrial fibrillation: the STROKESTOP study. *Circulation*.
430 2015;131(25):2176–84.
- 431 8. Hindricks G, Potpara T, Dagres N, Bax JJ, Boriani G, Dan GA, et al. 2020 ESC
432 Guidelines for the diagnosis and management of atrial fibrillation developed in collaboration
433 with the European Association for Cardio-Thoracic Surgery (EACTS). *Eur Heart J*.
434 2021;42(5):373–498.
- 435 9. Svennberg E, Engdahl J, Al-Khalili F, Friberg L, Frykman V, Rosenqvist M. Mass
436 screening for untreated atrial fibrillation: the STROKESTOP study. *Circulation*.
437 2015;131(25):2176–84.
- 438 10. Svennberg E, Stridh M, Engdahl J, Al-Khalili F, Friberg L, Frykman V, et al. Safe
439 automatic one-lead electrocardiogram analysis in screening for atrial fibrillation. *Europace*.
440 2017;19(9):1449–53.
- 441 11. Desteghe L, Raymaekers Z, Lutin M, Vijgen J, Dilling-Boer D, Koopman P, et al.
442 Performance of handheld electrocardiogram devices to detect atrial fibrillation in a cardiology
443 and geriatric ward setting. *EP Eur*. 2017 Jan 1;19(1):29–39.
- 444 12. Poulsen MB, Binici Z, Dominguez H, Soja AMB, Kruuse C, Hornnes AH, et al.
445 Performance of short ECG recordings twice daily to detect paroxysmal atrial fibrillation in
446 stroke and transient ischemic attack patients. *Int J Stroke*. 2017;12(2):192–6.
- 447 13. Wegner FK, Kochhäuser S, Ellermann C, Lange PS, Frommeyer G, Leitz P, et al.
448 Prospective blinded Evaluation of the smartphone-based AliveCor Kardia ECG monitor for
449 Atrial Fibrillation detection: The PEAK-AF study. *Eur J Intern Med*. 2020 Mar 1;73:72–5.
- 450 14. Pandiaraja M, Brimicombe J, Cowie M, Dymond A, Lindén HC, Lip GYH, et al.
451 Screening for atrial fibrillation: Improving efficiency of manual review of handheld
452 electrocardiograms. *Eng Proc*. 2020;2(1):78.
- 453 15. Viera AJ, Garrett JM. Understanding interobserver agreement: the kappa statistic. *Fam*
454 *Med*. 2005 May;37(5):360–3.

- 455 16. Mant J, Fitzmaurice DA, Hobbs FDR, Jowett S, Murray ET, Holder R, et al. Accuracy
456 of diagnosing atrial fibrillation on electrocardiogram by primary care practitioners and
457 interpretative diagnostic software: Analysis of data from screening for atrial fibrillation in the
458 elderly (SAFE) trial. *Br Med J*. 2007;335(7616):380–2.
- 459 17. Wegner FK, Kochhäuser S, Frommeyer G, Lange PS, Ellermann C, Leitz P, et al.
460 Prospective blinded evaluation of smartphone-based ECG for differentiation of
461 supraventricular tachycardia from inappropriate sinus tachycardia. *Clinical Research in*
462 *Cardiology*. 2021;110(6):905–12.
- 463 18. Kearley K, Selwood M, Bruel AV den, Thompson M, Mant D, Hobbs FR, et al. Triage
464 tests for identifying atrial fibrillation in primary care: a diagnostic accuracy study comparing
465 single-lead ECG and modified BP monitors. *BMJ Open*. 2014 Apr 1;4(5):e004565.
- 466 19. Lowres N, Neubeck L, Salkeld G, Krass I, McLachlan AJ, Redfern J, et al. Feasibility
467 and cost-effectiveness of stroke prevention through community screening for atrial
468 fibrillation using iPhone ECG in pharmacies: The SEARCH-AF study. *Thrombosis and*
469 *Haemostasis*. 2014;111(6):1167–76.
- 470 20. Koshy AN, Sajeev JK, Negishi K, Wong MC, Pham CB, Cooray SP, et al. Accuracy of
471 blinded clinician interpretation of single-lead smartphone electrocardiograms and a proposed
472 clinical workflow. *American Heart Journal*. 2018;205:149–53.
- 473 21. Racine HP, Strik M, Zande J van der, Alrub SA, Caillol T, Haïssaguerre M, et al. Role
474 of Coexisting ECG Anomalies in the Accuracy of Smartwatch ECG Detection of Atrial
475 Fibrillation. *Canadian Journal of Cardiology*. 2022 Nov 1;38(11):1709–12.
- 476 22. Pipilas DC, Khurshid S, Atlas SJ, Ashburner JM, Lipsanopoulos AT, Borowsky LH, et
477 al. Accuracy and variability of cardiologist interpretation of single lead electrocardiograms
478 for atrial fibrillation: The VITAL-AF trial. *American Heart Journal*. 2023 Nov 1;265:92–103.
- 479 23. Davidenko JM, Snyder LS. Causes of errors in the electrocardiographic diagnosis of
480 atrial fibrillation by physicians. *J Electrocardiol*. 2007 Sep 1;40(5):450–6.
- 481 24. Wong KC, Klimis H, Lowres N, Huben A von, Marschner S, Chow CK. Diagnostic
482 accuracy of handheld electrocardiogram devices in detecting atrial fibrillation in adults in
483 community versus hospital settings: a systematic review and meta-analysis. *Heart*. 2020 Aug
484 1;106(16):1211–7.
- 485 25. Ford C, Xie CX, Low A, Rajakariar K, Koshy AN, Sajeev JK, et al. Comparison of 2
486 Smart Watch Algorithms for Detection of Atrial Fibrillation and the Benefit of Clinician
487 Interpretation: SMART WARS Study. *JACC: Clinical Electrophysiology*. 2022 Jun
488 1;8(6):782–91.
- 489 26. Karregat EPM, Himmelreich JCL, Lucassen WAM, Busschers WB, van Weert HCPM,
490 Harskamp RE. Evaluation of general practitioners' single-lead electrocardiogram
491 interpretation skills: a case-vignette study. *Family practice*. 2021;38(2):70–5.
- 492 27. Himmelreich JCL, Karregat EPM, Lucassen WAM, Weert HCPM van, Groot JR de,
493 Handoko ML, et al. Diagnostic Accuracy of a Smartphone-Operated, Single-Lead
494 Electrocardiography Device for Detection of Rhythm and Conduction Abnormalities in
495 Primary Care. *The Annals of Family Medicine*. 2019 Sep 1;17(5):403–11.

496
497
498
499

500 **Figure Legends**

501 Figure 1: Data selection at the participant-level.

502 Figure 2: Data selection at the ECG-level.

503 Figure 3: Relationships between the level of agreement between cardiologists and factors at
504 the participant- and ECG-level. (* denotes a significant association)

505

Tables

Table 1: Agreement between cardiologists on participant-level AF diagnoses

		Cardiologist 2			
		AF	non-AF	cannot exclude AF	
Cardiologist 1	AF	44	26	2	72
	non-AF	5	78	4	87
	cannot exclude AF	1	26	4	31
		50	130	10	190

Table 2: Relationships between the level of agreement between cardiologists and factors at the participant and ECG-level

Factor	Categories	k_w	k	% <i>agree</i>	p-value
<i>Agreement at the participant-level</i>					
Age (years)	65-69	0.46 (0.15-0.73)	0.36 (0.09-0.62)	67.6	1.000
	70-74	0.47 (0.23-0.66)	0.42 (0.21-0.62)	66.0	
	75-79	0.46 (0.25-0.66)	0.42 (0.21-0.62)	66.0	
	80-84	0.42 (0.18-0.67)	0.40 (0.17-0.65)	66.7	
	85+	0.57 (0.30-0.80)	0.45 (0.21-0.72)	65.0	
Gender	Female	0.36 (0.18-0.54)	0.34 (0.17-0.51)	63.2	0.452
	Male	0.55 (0.40-0.67)	0.47 (0.34-0.59)	68.4	
Number of adequate-quality ECGs	0-54	0.21 (0.04-0.42)	0.21 (0.07-0.41)	47.7	0.001*
	55-66	0.33 (0.13-0.55)	0.26 (0.09-0.46)	56.9	
	67-108	0.74 (0.51-0.87)	0.64 (0.43-0.82)	80.4	
	109+	0.67 (0.46-0.84)	0.62 (0.40-0.80)	79.6	

Number of algorithm-identified possible AF ECGs	0-4	0.63 (0.32-0.85)	0.62 (0.34-0.86)	83.7	0.070
	5-9	0.20 (0.10-0.52)	0.19 (-0.07-0.52)	60.0	
	10-19	0.31 (0.06-0.53)	0.24 (0.06-0.45)	55.8	
	20-39	0.53 (0.32-0.74)	0.45 (0.25-0.66)	63.9	
	40+	0.47 (0.21-0.71)	0.38 (0.17-0.61)	66.7	
<i>Agreement at the individual ECG-level</i>					
Heart rate (bpm)	30-59	0.61 (0.49-0.71)	0.55 (0.44-0.66)	87.3	0.151
	60-69	0.63 (0.53-0.72)	0.57 (0.47-0.66)	86.5	
	70-79	0.62 (0.51-0.72)	0.54 (0.42-0.65)	88.5	
	80-89	0.44 (0.28-0.59)	0.38 (0.23-0.54)	84.2	
	90+	0.49 (0.36-0.61)	0.42 (0.30-0.55)	82.4	
RR-interval variability (%)	0.0-9.9	0.54 (0.44-0.64)	0.48 (0.39-0.58)	87.6	0.120
	10.0-19.9	0.57 (0.49-0.65)	0.48 (0.40-0.55)	84.3	
	20.0+	0.63 (0.53-0.71)	0.57 (0.48-0.66)	86.9	

ECG quality	Adequate-quality	0.61 (0.56-0.66)	0.55 (0.50-0.61)	88.0	0.000*
	Low-quality	0.17 (0.01-0.38)	0.13 (-0.02-0.31)	63.3	
Algorithm-identified possible AF?	No	0.24 (0.11-0.43)	0.23 (0.10-0.39)	92.5	0.000*
	Yes	0.59 (0.54-0.64)	0.53 (0.47-0.58)	82.8	

Table 3: Agreement between cardiologists on ECG-level AF diagnoses

		Cardiologist 2			
		AF	non-AF	cannot exclude AF	
Cardiologist 1	AF	144	84	7	235
	non-AF	28	1,424	37	1,489
	cannot exclude AF	7	93	19	119
		179	1,601	63	1,843

Figures

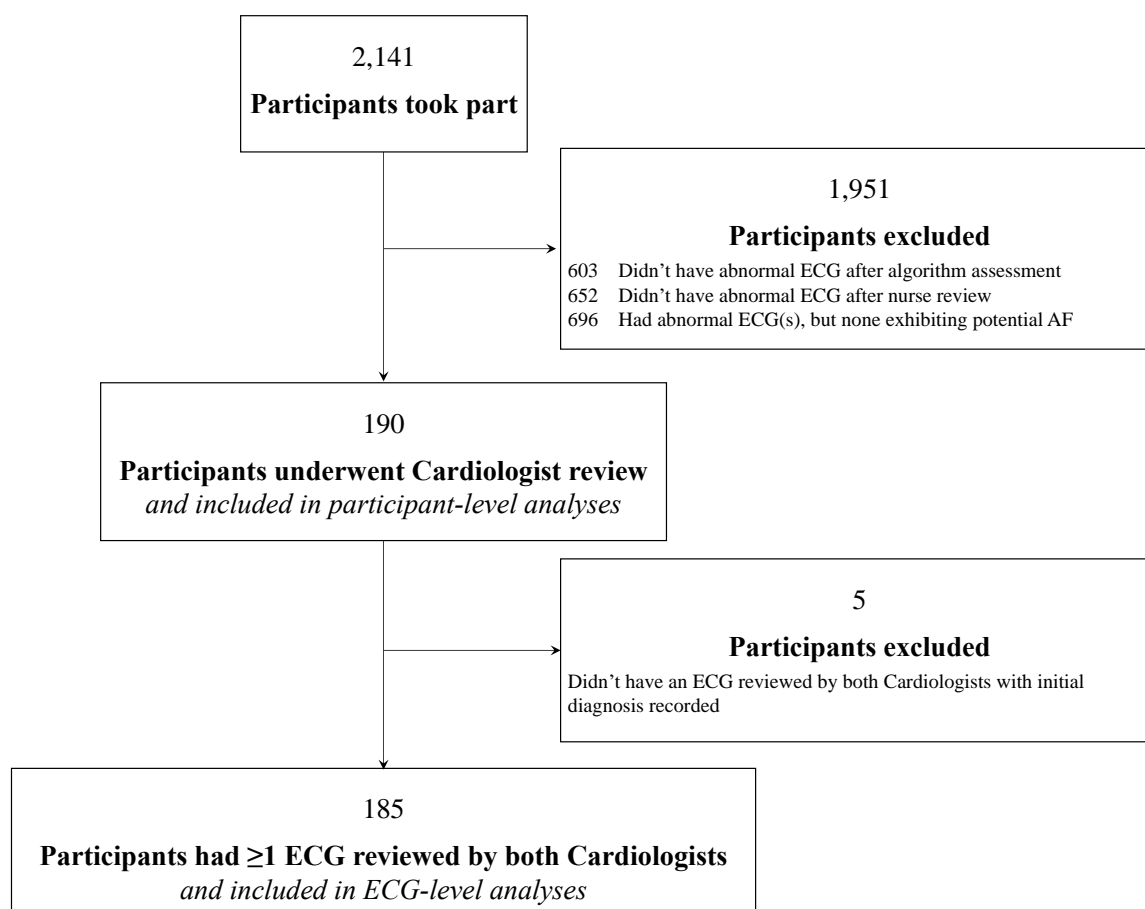
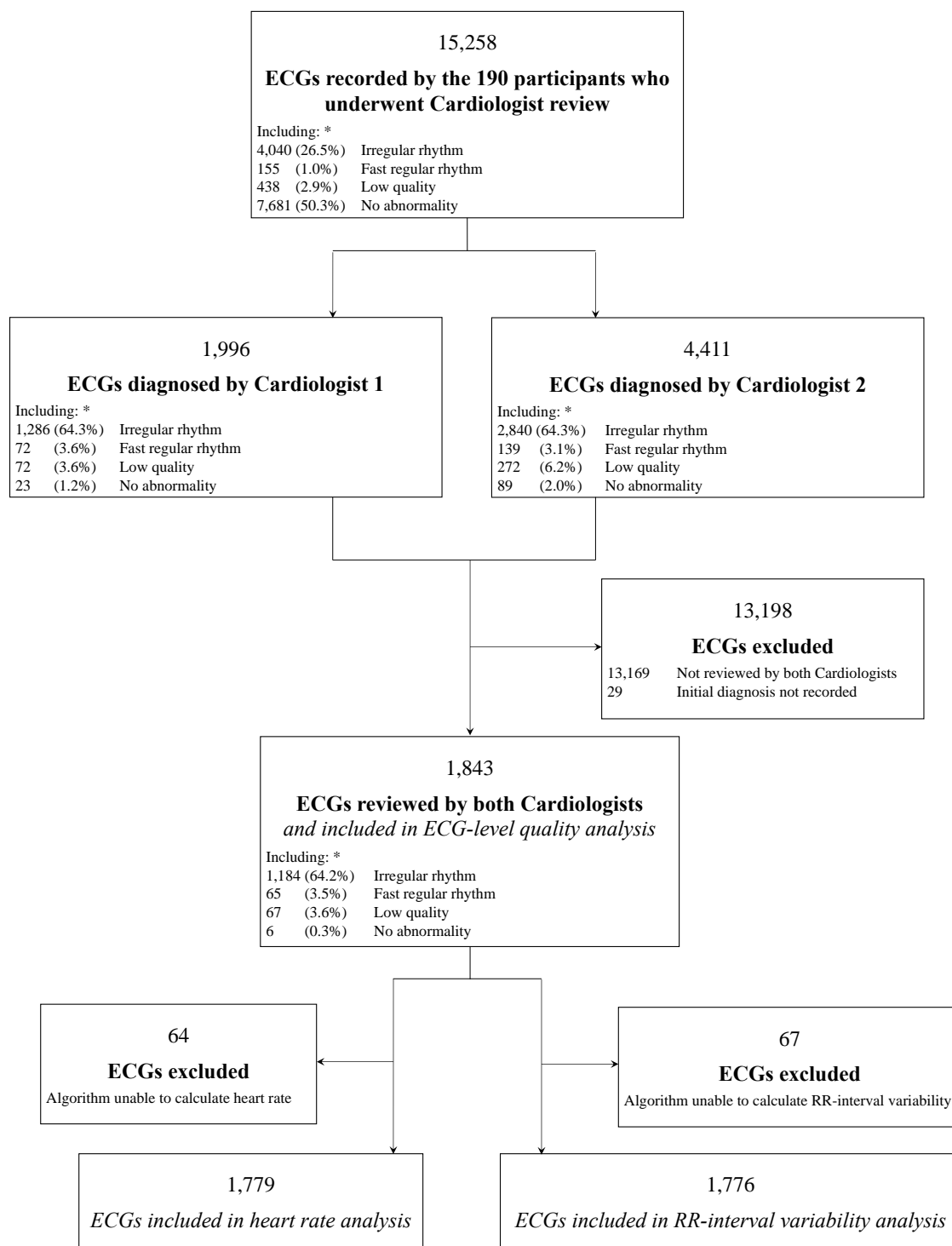


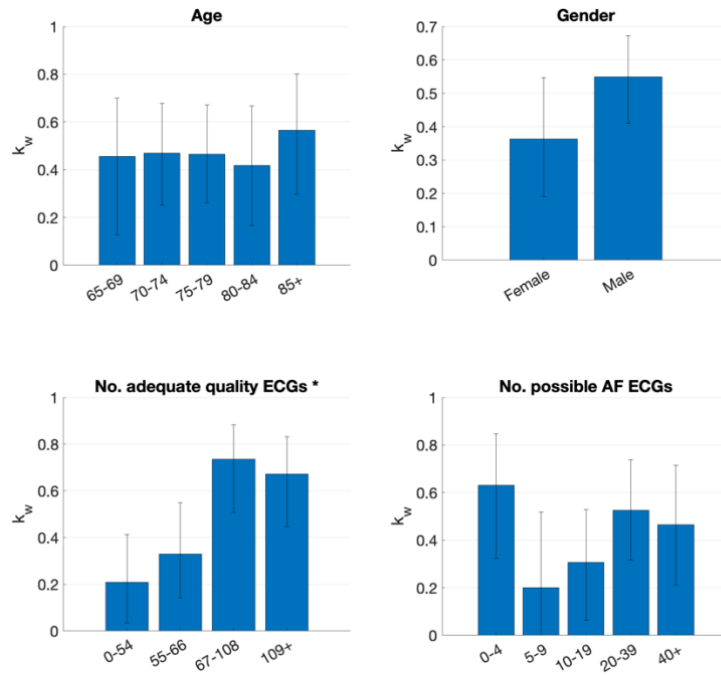
Figure 1: Data selection at the participant-level.



* ECGs could also be classified with other abnormalities by the computer algorithm, so these figures do not total 100%.
We considered 'Irregular rhythm' and 'Fast regular rhythm' as indicating signs of possible AF.

Figure 2: Data selection at the ECG-level.

Participant-level



ECG-level

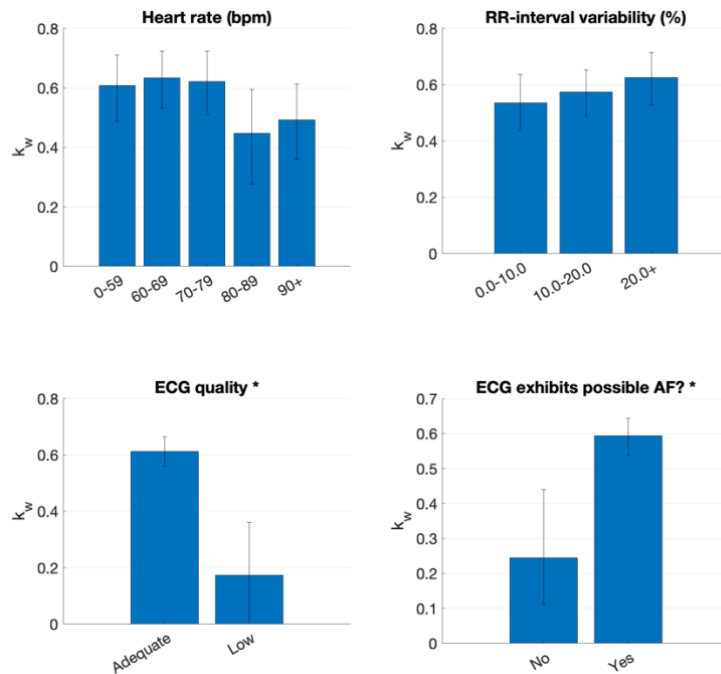


Figure 3: Relationships between the level of agreement between cardiologists and factors at the participant- and ECG-level. (* denotes a significant association)