# King's Research Portal

Document Version
Peer reviewed version

[Link to publication record in King's Research Portal](Link to publication record in King's Research Portal)

# Exploring Interpretability in Deep Learning Prediction of Successful Ablation Therapy for Atrial Fibrillation

1  **Shaheim Ogbomo-Harmitt[1], Marica Muffoletto[1], Aya Zeidan[1], Ahmed Qureshi[1], Andrew P.**
2  **King[1], Oleg Aslanidi[1]**

3  School of Biomedical Engineering and Imaging Sciences, King's College London, 3rd Floor Lambeth
4  Wing, St Thomas' Hospital, London, SE1 7EH, United Kingdom

5  **\*Correspondence:**
6  Oleg Aslanidi
7  oleg.aslanidi@kcl.ac.uk

10  **Abstract**

11  ***Background.*** Radiofrequency catheter ablation (RFCA) therapy is the first-line treatment for atrial
12  fibrillation (AF), the most common type of cardiac arrhythmia globally. However, the procedure
13  currently has low success rates in dealing with persistent AF, with a reoccurrence rate of ~50% post-
14  ablation. Therefore, deep learning (DL) has increasingly been applied to improve RFCA treatment for
15  AF. However, for a clinician to trust the prediction of a DL model, its decision process needs to be
16  interpretable and have biomedical relevance. ***Aim.*** This study explores interpretability in DL prediction
17  of successful RFCA therapy for AF and evaluates if pro-arrhythmogenic regions in the left atrium (LA)
18  were used in its decision process. ***Methods.*** AF and its termination by RFCA have been simulated in
19  MRI-derived 2D LA tissue models with segmented fibrotic regions (n = 187). Three ablation strategies
20  were applied for each LA model: pulmonary vein isolation (PVI), fibrosis-based ablation (FIBRO) and
21  a rotor-based ablation (ROTOR). The DL model was trained to predict the success of each RFCA
22  strategy for each LA model. Three feature attribution (FA) map methods were then used to investigate
23  interpretability of the DL model: GradCAM, Occlusions and LIME. ***Results.*** The developed DL model
24  had an AUC (area under the receiver operating characteristic curve) of $0.78 \pm 0.04$ for predicting the
25  success of the PVI strategy, $0.92 \pm 0.02$ for FIBRO and $0.77 \pm 0.02$ for ROTOR. GradCAM had the
26  highest percentage of informative regions in the FA maps (62% for FIBRO and 71% for ROTOR) that
27  coincided with the successful RFCA lesions known from the 2D LA simulations, but unseen by the
28  DL model. Moreover, GradCAM had the smallest coincidence of informative regions of the FA maps
29  with non-arrhythmogenic regions (25% for FIBRO and 27% for ROTOR). ***Conclusions.*** The most
30  informative regions of the FA maps coincided with pro-arrhythmogenic regions, suggesting that the
31  DL model leveraged structural features of MRI images to identify such regions and make its prediction.
32  In the future, this technique could provide a clinician with a trustworthy decision support tool.

33

## 1    Introduction

35  Atrial fibrillation (AF), the rapid, uncoordinated contraction of the atria, is a heart condition that affects
36  33 million people worldwide - making it the most common type of cardiac arrhythmia globally (Hart

37  and Halperin, 2001; Chugh et al., 2014). Currently, the precise mechanisms of AF are unclear.
38  However, there is evidence that ectopic electrical beats originating from the pulmonary veins (PVs)
39  can trigger AF (Chen et al., 1999). The triggers can then generate re-entrant drivers (rotors) that sustain
40  AF, and spatial fibrosis distributions in the left atria (LA) have been demonstrated to facilitate such
41  drivers (Morgan et al., 2016; Roy et al., 2020). A common treatment for AF is radiofrequency catheter
42  ablation (RFCA) therapy. RFCA involves using induced heat from a rapidly alternating current in a
43  catheter to ablate (isolate or destroy) the arrhythmogenic area of atrial tissue that harbours triggers or
44  rotors, thus restoring sinus rhythm and the mechanical function of the heart (Townsend and Sabiston,
45  2001). Presently, the success rate of RFCA is ~70% for paroxysmal AF - which is relatively high
46  (Oketani et al., 2012). However, the procedure is much less successful when dealing with persistent
47  AF, which has a reoccurrence rate of ~75% post-intervention. Therefore, with the high reoccurrence
48  rate of AF, there is a need for improvements within (Wang et al., 2017; Yubing et al., 2018).

49  Image-based computational modelling has been used to understand the structure-function relationship
50  that determines re-entrant atrial drivers for AF with the aim of improving RFCA outcomes. As a result,
51  computational methods have been introduced to improve RFCA outcomes, ultimately leading to the
52  FIRM (Focal Impulse and Rotor Modulation) mapping, which locates rotational activity around a
53  centre (rotor) from electroanatomical maps (Narayan et al., 2012a). The CONFIRM trial showed
54  patients that underwent FIRM-guided ablation maintained a higher freedom of AF (AF termination in
55  86% of patients) when compared to patients with conventional ablation strategy (AF termination in
56  20% of patients) (Narayan et al., 2012b). However, the multicentre REAFFIRM trial did not show
57  evidence that FIRM-guided ablation strategy is superior to pulmonary vein isolation (PVI) (Zhao et
58  al., 2019).

59  With the recent rise of artificial intelligence (AI), machine and deep learning (DL) have been applied
60  to patient medical imaging data and computational cardiac modelling with the aim to develop more
61  effective RFCA treatments. The applications of AI include predicting AF reoccurrence post-RFCA and
62  the origins of AF triggers and ablation (Kim et al., 2020; Liu et al., 2020; Firouznia et al., 2021; Roney
63  et al., 2022). Furthermore, Luongo et al. have applied machine learning to predict AF ablation targets,
64  but used 12-lead ECG data instead of medical imaging (Luongo et al., 2021). Other studies have also
65  leveraged the power of AI in AF by using DL with ECG data to estimate atrial fibrosis and to classify
66  AF from atrial flutter or tachycardia (Nagel et al., 2021; Rodrigo et al., 2022). Zololotarev et al. applied
67  AI to identify AF drivers from multi-electrode mapping, with the AI model then validated using optical
68  mapping; the results were comparable to the state-of-the-art with higher computational efficiency
69  (Zolotarev et al., 2020). Popescu et al. applied DL for arrhythmic sudden death prediction from clinical
70  and imaging data, which proved successful and achieved a concordance index of 0.83 and 0.74, and
71  10-year integrated Brier score of 0.12 and 0.14, respectively (Popescu et al., 2022).

72  However, DL is limited by its black-box nature. This is an issue when considering the European
73  Union's General Data Protection Regulation (GDPR), as any algorithmic decision used in patient care
74  requires an explanation for transparency (Mourby et al., 2021). Moreover, clinicians have also argued
75  that if AI can outperform human diagnosis, understanding the AI model's decision process will be
76  beneficial in discovering new biological processes and furthering medical knowledge (Watson et al.,
77  2019). Moreover, it will increase confidence in the AI-generated results, which means the clinicians
78  are more likely to trust and leverage them. Hence, this has led to the growing field of deep learning
79  interpretability for medical imaging analysis, where methods such as concept learning models, latent
80  space interpretation and attribution maps have been applied to many medical fields (Salahuddin et al.,
81  2022). Organisations have also expressed an interest in AI interpretability, e.g., the Avicenna Alliance
82  (AA) and the Virtual Physiological Human Institute (VPHI). The AA and VPHI aims are to promote

83  the synergy of AI and in-silico modelling into healthcare, while providing policy makers and regulators
84  with directions towards applying these technologies safely in clinics, including AI interpretability
85  (Liesbet Geris et al., n.d.).

86  Muffoletto et al. were the first to apply DL to directly informing a clinician to treat AF using RFCA
87  therapy and developed a convolutional neural network (CNN) to predict suitable in-silico ablation
88  strategies for a given patient, using synthetic tissue-based atrial models with randomly distributed
89  fibrotic patches. The approach proved effective (79% accuracy) and illustrated the proof-of-concept
90  (Muffoletto et al., 2019). Ultimately, this led to the approach being applied to MRI-derived data to
91  predict the patient-specific optimal RFCA strategy. As a result, the developed CNN had a 100%
92  accuracy for classifying optimal fibrosis- (FIBRO) and rotor-based (ROTOR) strategies success and
93  33% accuracy for the PVI strategy (Muffoletto et al., 2021).

94  Currently, research into interpretability for DL-based AF management is very limited. For example,
95  one study by Alhusseini et al. used gradient-weighted class activation mapping (GradCAM) to show
96  that their feature attribution (FA) map closely replicated rules used by clinicians. However, only one
97  method was validated within this study, and a comparison with other methods was not investigated.
98  Furthermore, the study used spatial maps of the activation phase derived from electrocardiogram data
99  from a basket catheter. Hence, there has been no investigation into DL interpretability for models which
100 use medical imaging data to make explainable predictions for cardiac arrhythmias and anti-arrhythmic
101 treatments (Alhusseini et al., 2020).

102 In this study, we present a novel qualitative and quantitative comparison of established DL
103 interpretability methods for medical imaging and image-based cardiac modelling of RFCA, as well as
104 new quantitative metrics to assess interpretability of FA maps for the image-based cardiac models.

105

## 106  2     Methods

### 107  2.1   Overview

108 We propose a DL approach to 1) accurately predict the outcomes of RFCA therapy based on image-
109 based modelling and simulations and 2) interpret the decision process of the DL model. To achieve
110 this, standardised 2D LA models with patient-specific distributions of fibrosis were derived from late
111 gadolinium-enhanced (LGE) MR imaging data. Simulations of AF and its termination with three RFCA
112 strategies were performed, the DL model was applied to predict the success of each strategy, and the
113 RFCA simulation results were compared with DL interpretability maps to identify proarrhythmogenic
114 locations. Three established interpretability approaches were also compared qualitatively and
115 quantitatively to interpret the DL model's predictions.

### 116  2.2   Data Acquisition and Pre-processing

117 The datasets used in this study were derived from 122 LGE MRI patient scans: 86 datasets with spatial
118 resolution of 0.625x0.625x0.625 $mm^3$ were acquired from the Atrial Segmentation Challenge at the
119 STACOM 2018 workshop (Xiong et al., 2021); additionally, 36 LGE MRI images were collected at
120 St. Thomas' Hospital London with resolution of 1.3x1.3x4 $mm^3$ (specifically, 18 AF patients were
121 scanned both pre-and post-intervention) (Chubb et al., 2018).

122 Generating 2D LA models with fibrosis first required manual segmentation of patient LGE MRI data
123 to produce 3D patient-specific endocardial LA surface meshes. The LGE MRI image intensities were
124 then mapped to these models and the image intensity ratio thresholding technique was applied to
125 quantify and visualise LA fibrosis (Roy et al., 2020). Finally, the 3D LA fibrosis maps were unwrapped
126 using the LA standardised unfold mapping technique to produce models in the 2D LA disk format for
127 use as input to the DL network, as shown in Figure 1A (Williams et al., 2017; Qureshi et al., 2020).

128 Furthermore, to increase the size of the dataset, synthetic 2D LA disks were generated by weighted-
129 averaging of the patient datasets to vary the fibrosis distribution and PVs. The creation of synthetic
130 disks consisted of three steps. First, 65 MRI images were extracted from the STACOM 2018 dataset
131 and were each weighted by assigning a random weight (between 0 to 1) to all voxels of a given image;
132 the weighted-average of all images was thresholded (Case xA in Figure 1B). This number was chosen
133 as less than 65 would result in low variability in the synthetic tissues and more than 65 would result in
134 most of the synthetic tissues being covered in fibrosis. Supplementary Figure S1 illustrates that
135 selecting the 65 LA tissues in generating the synthetic LA tissues would result in a mean fibrotic tissue
136 percentage of approximately 50%. Thus, 65 corresponds to a folding point of this sigmoidal
137 dependence, and any number above 65 would lead to a majority of tissue being fibrotic. Then the
138 extracted fibrosis distribution was further augmented by applying one or multiple affine
139 transformations (translation, rotation and flipping). The fibrosis threshold value and the types of
140 transformation were randomly selected. Lastly, the PVs were varied by assigning one of 6 different
141 variants, which included changing PV size and position (Case xB in Figure 1B) (Muffoletto et al.,
142 2021). This resulted in a total of 199 synthetic 2D LA tissue models in addition to the 122 patient-
143 specific models, totalling 321 2D LA tissue models.

144 **2.3 Atrial Tissue Modelling and AF simulation**

145 Equation (1) represents the Fenton-Karma semi-physiological model, which consists of three ionic
146 currents representing the overall ion current in the electrical dynamics of atria cells; $I_{fi}$ represents the
147 fast inward current $Na^+$, $I_{so}$ is the slow outward current $K^+$ and $I_{si}$ is the slow inward current $Ca^+$
148 (Fenton and Karma, 1998):

149
$$I_{ion} = I_{fi} + I_{so} + I_{si} \tag{1}$$

150 Equation (2) is the standard monodomain equation to describe electrical wave propagation.

151
$$\frac{\partial V_m}{\partial t} = \nabla . D \nabla V_m - \frac{I_{ion}}{C_m} \tag{2}$$

152 $V_m$ is the membrane potential, $C_m$ is the membrane capacitance, $D$ is a tensor that represents the
153 diffusion of the electrical coupling within tissue. Equation (2) with ion current determined in equation
154 (1) was solved using the forward Euler method with a finite-difference approximation of the Laplacian.
155 Therefore, equation (1) and equation (2) were solved using each 2D tissue disk as a spatial domain to
156 simulate electrical waves sustaining AF. Such waves in the form of rotors were generated using the
157 standard cross-field protocol at 28 $ms$ into the simulation (Tobón et al., 2014). The numerical
158 integration steps were 0.01 ms time step and 0.3 $mm$ spatial step. Additionally, healthy tissue had a $D$
159 value of 0.1 $mm^2 s^{-1}$ to match the physiological value of healthy myocardium tissue. Fibrotic tissue
160 had $D$ value of 0.015 $mm^2 s^{-1}$.

161 The three ablation strategies were simulated to terminate persistent AF: PVI, FIBRO and ROTOR
162 strategies; details of the simulations have been published previously (Muffoletto et al., 2021). The
163 FIBRO strategy involved ablating the perimeter of the fibrotic tissue, while PVI consisted of ablating
164 the circumference of the PVs and ROTOR ablated the phase singularities of the electrical wave. The
165 ablation strategy was deemed successful for a tissue if AF was terminated within 2000 *ms* and less than
166 40% of the tissue was ablated (Muffoletto et al., 2021). Therefore, using the stated simulation pipeline,
167 the success of the three RFCA strategies was determined for AF simulations in the 2D LA tissues
168 (including patient MRI derived and synthetic data). Furthermore, since multiple strategies can be
169 successful/unsuccessful for a given 2D LA tissue, the classification task was multi-label.

170 **2.4   Deep Learning**

171 We employed the CNN with hyperparameters (parameters and number of convolutional and fully
172 connected layers) based on the study by Muffoletto et al. as the basis of our interpretability framework
173 (Muffoletto et al., 2021). The hyperparameters were tuned by Muffoletto et al. by performing 24
174 experiments which involved changing number of layers, filter size of convolutional layers and dropout
175 rate. The optimal hyperparameters were chosen by selecting the DL model with the highest average
176 accuracy across a 5-fold cross-validation. The CNN consisted of four convolutional layers of 32x32
177 filters, each followed by Rectified Linear Unit (ReLU) activation and max pooling with a pool size of
178 two. These are followed by three linear layers (2048, 128 and 3 units, respectively) and another ReLU
179 activation. A Dropout layer followed this at a rate of 0.8 and a sigmoid function (Paszke et al., 2019).
180 Since we address a multi-label classification problem (i.e., multiple ablation strategies), we modified
181 the loss function to be a mean-squared error tailored to perform multi-label classification for the three
182 ablation strategies (Figure 1).

$$MSE(\boldsymbol{y}_{score}, \boldsymbol{y}) = \sum_{i=0}^{N} \frac{\boldsymbol{y}_{score}^{i} - \boldsymbol{y}^{i}}{N} \tag{3}$$

184 Equation (3) is the mean-squared error function formulation, where $\boldsymbol{y}_{score}$ is the predicted class score
185 array and $\boldsymbol{y}$ is the RFCA strategy success ground truth (where 1 = success and 0 = unsuccessful). Here,
186 $\boldsymbol{N}$ represents the number of classes/strategies (three in this study) and $\boldsymbol{i}$ is the index of a class in the
187 class score array. To train and effectively test the CNN, a leave-one-out cross-validation was used
188 where the total dataset was split into two sets: a hold-out test set and training set. The training set was
189 then split into five folds, where four folds were used to train the CNN, and the last fold was used as a
190 validation set to select the optimal CNN model state (i.e. the model with the lowest loss during training)
191 (Raschka, 2018; Muffoletto et al., 2021). In total, there were 271 2D LA tissues in the leave-one-out
192 cross-validation dataset (96 MRI derived and 175 synthetic). Within each fold the DL model was
193 trained for 100 epochs using an ADAM optimiser with a learning rate of 1e-4 (Kingma and Ba, 2014).
194 For each fold, the optimal model was tested on the hold-out test set of 50 2D LA tissues (26 MRI
195 derived and 24 synthetic) from the total dataset to evaluate the DL model's performance. Pre- and post-
196 ablation images were not split during cross-validation, as there was little similarity between the
197 respective fibrosis distributions (see Supplementary Materials Section 2 and Supplementary Figure
198 S2).

199 **2.5   Interpretability**

200 Three popular local post-hoc interpretability methods were used to interpret the CNN's predictions -
201 GradCAM, occlusions and local interpretable model-agnostic explanations (LIME) (Zeiler and Fergus,
202 2014; Ribeiro et al., 2016; Selvaraju et al., 2017; Kokhlikyan et al., 2020). GradCAM and LIME were

203    chosen as they are widely used saliency maps in DL medical image analysis (Magesh et al., 2020;
204    Graziani et al., 2021; Patel et al., 2021; Mahapatra et al., 2022), while occlusions is one of the first
205    saliency map methods used for DL computer vision. Each method evaluates feature attribution using
206    different approaches: GradCAM uses gradient information, LIME uses an interpretable model within
207    a local space and the occlusions method uses perturbations.

208    The DL model state from the most accurate fold of the leave-one-out cross-validation was used to
209    produce the FA maps for the three methods on the hold-out test set. The GradCAM method was applied
210    to the last convolutional layer of the CNN. Each FA map was thresholded above the respective map's
211    average FA to highlight the most informative features. Three metrics were evaluated to quantitatively
212    analyse the informative regions of each FA map: Jacquard index (IoU), lesion percentage and non-
213    arrhythmogenic tissue (NAT) percentage. The IoU was evaluated by calculating IoU of the informative
214    regions of a FA map and lesions of a given ablation strategy. Lesion percentage was evaluated by
215    calculating the percentage of lesions of a given ablation strategy within the informative regions.

216    The motivation for analysing the lesion percentage was to determine if the DL model focused on
217    clinically relevant features. The number of the lesions (unseen by the DL model but known from
218    simulations – and known to clinicians when ablating a patient) found in a FA map's informative region
219    is a relevant metric, as such lesions are associated with arrhythmogenic regions in atrial tissue. Thus,
220    PVI lesions isolate the area of the initial arrhythmogenic triggers, FIBRO lesions aim to isolate the
221    fibrotic tissue border where AF reentrant drivers commonly reside, and ROTOR lesions directly target
222    such reentrant drivers. Therefore, the ability of DL model to predict lesion locations (again, without
223    seeing such lesions during training) should help the clinician to understand and trust these predictions.

224    Lastly, the NAT percentage was calculated by finding the percentage of healthy tissue (with no lesions
225    or fibrosis) within the informative regions of a FA map. NAT percentage was evaluated to assess how
226    much of the clinically irrelevant features were highlighted as informative by the DL model.

227

228    **2.5.1 GradCAM**

229    GradCAM uses the gradient from a given convolutional layer to measure FA for a particular decision
230    of interest. GradCAM is an improvement of the class activation map (CAM) method. CAM produces
231    a localisation map for an image classification model, utilising a specific architecture where globally
232    averaged pooled convolutional feature maps are fed directly into a softmax layer. GradCAM improves
233    on CAM by being architecture-independent, and it can be applied to any CNN. Furthermore, a study
234    by Adebayo et al. implemented a sanity check of GradCAM through a model parameter and data
235    randomisation test. It demonstrated that GradCAM's saliency maps could support tasks that require
236    explanations that are faithful to the model and the data generation process (Adebayo et al., 2018).

$$\alpha^c = \frac{1}{z} \sum_i \sum_j \frac{\partial y^c}{\partial A_{ij}} \qquad (4)$$

237

238    Feature attribution, $\alpha_{ij}^c$ ($i$ and $j$ are the indices of the feature in a FA map), of a given class $c$ is calculated
239    in GradCAM by evaluating the partial derivative of the score of class $c$ and a feature from activation
240    map of a given convolutional layer $A_{ij}$. The result of evaluating the partial differential for each feature
241    is then pooled globally by dividing each element of the FA map by the total number of features to find
242    the final FA map (Selvaraju et al., 2017).

243 **2.5.2 LIME**

244 The core idea of LIME is to explain predictions of any classifier faithfully by learning an interpretable
245 model locally around the prediction. LIME achieves this by generating simulated data points around
246 an instance through random perturbation and weighting them as a function of proximity to the original
247 data points, fitting a sparse linear model to the predicted responses from the perturbed points and using
248 the sparse linear model as an explanation model (i.e., weights of features in linear model).

$$\xi(x) = \underset{g \in G}{\text{argmin}} \, \mathcal{L}(f, g, \pi_x) + \Omega(g) \tag{5}$$

$$\mathcal{L}(f, g, \pi_x) = \sum_{z, z' \in Z} \pi_x(z)(f(z) - g(z'))^2 \tag{6}$$

251 **FA** $\xi(x)$ of given features $x$ is calculated in LIME by minimising the loss function $\mathcal{L}$ and complexity,
252 $\Omega(g)$, of the function $g$ (a model from a class of possibly interpretable models). In essence $\mathcal{L}$ is a
253 function that measures how unfaithful the function $g$ is at approximating $f$ (the model being explained)
254 in the local space defined by $\pi_x$. Equation (6) shows how the loss function uses the L2 distance to
255 measure how unfaithful function $g$ is at approximating $f$, where $z$ is sample from $x$, $z$ is the set
256 perturbed samples of $x$ with associated labels and $z'$ is perturbed sample from set $z$ (Kokhlikyan et al.,
257 2020).

258 **2.5.3 Occlusions**

259 Occlusions is a perturbation-based approach to calculate FA, which involves perturbing the feature
260 space with a rectangular region and evaluating the difference of class score from a given class
261 prediction by the perturbation. FA is then assigned by looking at the feature in the multiple rectangular
262 regions it is in and averaging the multiple class score differences (Ancona et al., 2017). The occlusion
263 FA method was based on an occlusion sensitivity analysis used to validate a DL interpretability method
264 by Zeiler et al. (Zeiler and Fergus, 2014).
265

266 **3    Results**

267 **3.1.1 Dataset Analysis**

268 In the dataset comprising of 122 LA tissues derived from MRI data, the PVI strategy led to successful
269 AF termination in only 11.6% of cases, while 88.4% resulted in failed terminations. Meanwhile, the
270 FIBRO and ROTOR strategies resulted in 42.6% and 74.4% successful terminations, respectively.
271 Notably, FIBRO demonstrated the most balanced AF termination outcomes, whereas ROTOR and PVI
272 exhibited a similar level of misbalance in the outcomes. In the larger dataset consisting of 321 LA
273 tissues, including both MRI-derived and synthetic data, the PVI strategy achieved successful AF
274 termination in 27.1% of cases, demonstration a positive impact of augmentation. The FIBRO and
275 ROTOR strategies also resulted in 58.3% and 75.7% successful terminations, respectively.

276 **3.1.2 Convolutional Neural Network Performance**

277 The success of the FIBRO ablation strategy was predicted most accurately by the CNN, as shown in
278 Table 4, where the FIBRO class has the highest AUC score and the most balanced recall and precision
279 scores. Furthermore, the FIBRO strategy also had the highest AUC score when predicting ablation
280 success exclusively on the real data (Table 2). PVI had the second-highest AUC score on mixed real
281 and synthetic data, as well as exclusively on real data. Meanwhile, ROTOR had a comparable AUC

282  score to PVI on the real and synthetic data but performed worse exclusively on the MRI-derived data
283  (Table 2).

284  However, the CNN struggled to predict successful AF termination cases by PVI, which is reflected in
285  the low recall and F1 score in Table 1. Even though there was a similar class imbalance in ROTOR
286  compared to PVI, the CNN was able to predict the successful and failed AF termination cases to a
287  reasonable degree (see recall and F1 score in Table 1). Lastly, the CNN had a significantly higher AUC
288  score ($p < 0.05$) when trained and predicted on a dataset comprised of synthetic and MRI derived data
289  compared to training exclusively on MRI derived data (Table 3). This was confirmed using a one-sided
290  t-test (PVI: $p = 0.030$; FIBRO: $p = 3.5e-05$; ROTOR: $p = 6.15e-06$). This was due to the increased
291  dataset size when combining the real and synthetic data as the CNN has more training examples –
292  effectively improving the task's generalisation. Notably, incorporating synthetic data has improved
293  accuracy in predicting the outcomes of PVI. When trained exclusively on MRI-derived data, the model
294  showed a zero F1-score for PVI, attributed to significant class imbalance. This resulted in the model
295  predicting unsuccessful AF termination for all PVI cases, explaining the precision score of 1.0.
296  However, integrating synthetic data into the dataset improved the model's ability to classify successful
297  ablation for PVI (F1 score of $0.42 \pm 0.06$), due to the 15.5% increase in successful PVI cases in the
298  dataset. This allowed the model to improve its classification of successful AF termination by PVI.

### 3.1.3 Qualitative Interpretability Analysis

300  As shown in Table 4, GradCAM was characterised by the highest lesion percentage and IoU metrics
301  for the FIBRO and ROTOR strategies. Additionally, Figure 3 shows that in FA maps obtained with
302  GradCAM for ROTOR and FIBRO, the informative regions coincided with most ablation lesions.
303  Figure 3 also illustrates that GradCAM had the lowest NAT percentage for the FIBRO and ROTOR
304  strategies, as the FA maps did not highlight large, but clinically irrelevant regions of healthy tissue –
305  whereas LIME and occlusions did. For the PVI strategy, the occlusions method provided FA maps
306  with the greatest lesion percentage, and LIME provided FA maps with the highest IoU score.

### 3.1.4 Quantitative Interpretability Analysis

309  Using the Wilcoxon signed-rank test, the ROTOR strategy lesion percentage for GradCAM was
310  significantly greater ($p < 0.017$ using Bonferroni correction) than that for occlusions, but not for LIME
311  ($p = 3.1e-8$ and $p = 0.0253$, respectively). Moreover, for the FIBRO strategy, the lesion percentage for
312  GradCAM was significantly higher than that for the occlusions method, but again not for LIME ($p =$
313  $4.0 e-6$, $p = 0.06$, respectively). However, the IoU scores for GradCAM were significantly greater ($p <$
314  $0.017$) than those for occlusions and LIME for ROTOR ($p = 3.3e-6$ and $p = 2.1e-9$, respectively) and
315  FIBRO ($p = 4.2e-6$ and $p = 1.6e-9$, respectively). GradCAM also had a significantly less NAT
316  percentage ($p < 0.017$) than occlusions and LIME for ROTOR ($p = 5.5e-05$ and $p = 2.3e-09$,
317  respectively) and FIBRO ($p = 1.2 e-5$ and $2.3e-6$, respectively).

318  Therefore, GradCAM produced more interpretable FA maps than LIME (for FIBRO and ROTOR) as
319  the informative regions were more focused on areas with a high number of ablation lesions – reflected
320  in GradCAM having a significantly greater IoU score than LIME (Figures 5 and 4). Furthermore,
321  GradCAM was also more interpretable in a sense that its FA maps highlighted less regions that were
322  non-arrhythmogenic, and hence it had a significantly less NAT percentage than LIME and occlusions
323  (Figure 6).

324  For the PVI strategy, the occlusions method provided FA maps with the greatest lesion percentage and
325  LIME FA maps had the highest IoU score. The difference in best FA map methods in terms of lesion

326   percentage and IoU score can be seen in Table 4, as informative regions in the occlusions' FA maps
327   cover a vast area highlighting the ablation lesions but are not local to the PVs. Meanwhile, the LIME
328   FA map highlights areas around the PVs, but does not cover many ablation lesions.

330   Supplementary Figures S3, S4 and S5 show the difference in the mean score of each interpretability
331   metric for correct and incorrect classifications of AF termination for each ablation strategy and FA
332   method on the hold-out test set. This analysis shows no clear or consistent relationship between
333   interpretability and model accuracy. The mean interpretability scores reflect this, as they were similar
334   across the correct and incorrect classification groups. Additionally, the mean interpretability score
335   variability is inconsistent across each ablation strategy FA method and interpretability metric - further
336   illustrating no relationship between interpretability and accuracy.

### 3.1.5 Feature Attribution Thresholding Sensitivity Analysis

339   The findings presented above show little dependence on the threshold between informative and
340   uninformative regions. As shown in Figure 9, when the threshold value is set to 25% above and below
341   the average feature attribution, Grad-CAM still has the highest lesion percentage and IoU compared to
342   LIME and Occlusions for the ROTOR and FIBRO strategies. GradCAM still had a lower NAT
343   percentage for FIBRO and ROTOR when the threshold value was 25% below the average FA.
344   However, occlusions had a lower NAT percentage for FIBRO and ROTOR when the threshold value
345   was above 25% of the average FA. Occlusions had a lower lesion percentage and IoU, which shows
346   that GradCAM was more interpretable when the threshold was 25% above the average FA.

### 3.1.6 Population-level Interpretability Analysis

349   Figure 10 compares the average GradCAM FA maps for ROTOR, FIBRO and PVI with the average
350   fibrosis density across the 2D LA tissue disks. It shows that the high FA regions in the average FA
351   map for ROTOR (Figure 10B) and FIBRO (Figure 10C) correspond with dense fibrotic areas (Figure
352   10A). Furthermore, there was a similar good correspondence between the average GradCAM FA maps
353   for ROTOR and FIBRO (Figure 11B and 11D) and the respective average lesions across the 2D LA
354   tissue disks (Figures 11A and 11C). Unsurprisingly, the average GradCAM FA map for PVI (Figure
355   10D) showed relatively small correspondence to areas with high fibrosis density areas.

### 4       Discussion and Conclusion

358   Predicting RFCA outcomes from imaging data is a challenging task, as shown by Kim et al., who
359   predicted AF recurrence post-RFCA with a 0.61 accuracy from a CNN which used a combination of
360   MRI data and patient demographics (Kim et al., 2020). Moreover, Roney et al. applied machine
361   learning to predict in-silico AF recurrence after multiple ablation strategies (Roney et al., 2018, 2020).

362   Therefore, developing a successful DL model to predict RFCA outcomes in AF simulations is the
363   natural first step to predict real RFCA outcomes in AF patients. Hence, this study (i) demonstrates a
364   multi-label classification CNN for the success of ablation strategies in patient-specific simulations of
365   AF, with AUC scores of $0.92 \pm 0.02$ for FIBRO, $0.78 \pm 0.04$ for PVI and $0.77 \pm 0.02$ for ROTOR, and
366   (iii) explores different methods of DL interpretability in the classification, with GradCAM shown to

367 provide the most interpretable FA maps for the ROTOR and FIBRO strategy, suggesting that the DL
368 model utilises pro-arrhythmogenic regions to make its prediction. This is further supported by the
369 population-level interpretability analysis, as average FA maps for ROTOR and FIBRO are focused on
370 areas with high fibrotic density. This can be explained by the fact that the respective ablation lesions
371 are primarily located within these areas. Hence, the DL model can learn to predict AF termination
372 outcomes by implicitly leveraging pro-arrhythmogenic regions related to a given strategy. Importantly,
373 locations of the ablation lesions have not been explicitly used in the CNN's learning process.

374 It is worth noting that classification of the PVI strategy was difficult to interpret. A possible reason for
375 this difficulty is that the PVI strategy in the clinic is based on ablating PV triggers that typically initiate
376 AF. However, these initial PV triggers were not present in the 2D LA tissue models. Therefore, the
377 three FA methods could not produce interpretable maps in this case.

378 A possible explanation for why GradCAM performed better than the other methods is that LIME is
379 susceptible to unstable generated interpretations due to random perturbations and feature selection.
380 Moreover, LIME and occlusions are not class discriminative – meaning that they cannot localise the
381 class (RFCA strategy) within the feature space. GradCAM is gradient-based (does not randomise
382 parameters to obtain FA maps) and is class discriminative, allowing it to localise pro-arrhythmogenic
383 regions more faithfully than LIME and occlusions (Selvaraju et al., 2017; Zafar and Khan, 2021).

384 The RFCA strategy that has the highest magnitude of lesion percentage and lowest magnitude of NAT
385 percentage (ROTOR) also had the lowest AUC score in testing (Table 1), showing that the
386 interpretability of a FA map does not increase with the accuracy of the strategy's prediction. This
387 observation demonstrates that the need for interpretability in RFCA strategy prediction likely goes
388 beyond FA, and in future work, we will investigate the incorporation of confidence in prediction
389 outputs to enable our method to be used as a decision support tool to help clinicians select the
390 appropriate therapy. Since Varela et al. showed that LA anatomy is a significant factor in prediction of
391 AF recurrence post ablation (Varela et al., 2017a), the DL approach of the study should be extended to
392 3D LA images and simulations. Future work should also focus on using exclusively real patient LA
393 data and investigating intrinsically interpretable DL models such as ICAM (Bass et al., 2022).

394 Note that 2D LA disks were used in this study due to the efficiency in providing the needed proof of
395 concept and had clear advantages over extremely computationally-intensive 3D atrial simulations.
396 Moreover, the standardised 2D unfolded LA images allowed for generation of a large number of
397 additional synthetic images, which is crucial for training CNNs. Hence, image-based 2D LA models
398 provided a sensible balance between realistic details (such as fibrosis distributions) and computational
399 efficiency (i.e., the ability to run a large number of simulations and train the CNN). Previous work has
400 shown that atrial wall thickness is distributed more or less evenly in the LA outside of PVs and that
401 slow conduction in fibrotic areas is the main determinant of the rotor dynamics (Varela et al., 2017b;
402 Roy et al., 2018).

403 Another worthwhile direction is applying an approach based on counterfactual explanations, which
404 alters the input's feature space to change the classifier's prediction. Mertes et al. has applied this
405 approach to a generative adversarial network and showed its superiority to LIME in an X-ray imaging
406 study of pneumonia (Mertes et al., 2022). This research utilised over 100 non-medical experts for the
407 evaluation, which ultimately should become a standard for any interpretability study.

408 Our original approach to the evaluation is based on using a large number of 2D LA tissue models with
409 tractable features (rather than a large number of experts) to understand the predictions of the DL model.

410    Simulations of the test set of 50 2D LA tissue models reveal the important features determining the
411    success of each given RFCA strategy, such as the precise locations of ablation lesions and underlying
412    structural features. This evaluation shows that GradCAM best characterises if a DL model leverages
413    relevant features in its predictions. The fact that GradCAM highlights relevant features and does not
414    highlight healthy tissue devoid of such features is illustrated in Figures 3, 7 and 8 – and supported by
415    numerical metrics calculated using all 50 LA tissue models and summarised in Table 4.

416    The EU's GDPR requires an explanation for any algorithmic decision used in patient care; we believe
417    our work represents a significant step to meet this requirement. Most of the ablation lesions in our
418    study coincided with informative regions of the GradCAM FA maps (specifically, for ROTOR and
419    FIBRO, see Figures 8 and 7), whereas healthy, non-arrhythmogenic tissue (NAT) was outside of these
420    informative regions. This suggests that the DL model can learn from structural features of patient MR
421    images even without knowledge of the LA function. The explanation is that the structural features
422    constitute pro-arrhythmogenic LA regions (e.g., fibrotic regions are well-known for their ability to
423    harbour rotors sustaining AF) that need to be targeted by ablation. Such mechanistic explanations
424    should increase clinician's confidence in using the DL predictions in future.

425    This study's analysis also suggests that there is no clear relationship between a model's interpretability
426    and accuracy, which opens future directions of research into the relationship and interaction between
427    a model's performance and explainability. Another interesting investigation would be into how FA
428    maps can be used as model feedback to improve its performance. To our knowledge, no study has
429    investigated the application of interpretability feedback for DL model design and development for
430    biomedical applications. Bell et al. investigated the trade-off between accuracy and explainability for
431    black box and interpretable models. They showed that the trade-off is inconsistent, and in some cases
432    models with high explainability can also have high accuracy - but in others higher explainability comes
433    at the expense of low accuracy (Bell et al., 2022).

434    Importantly, the purpose of FA maps is not to be directly applied in the clinic to predict ablation lesions
435    in a patient – but to explain why the DL approach is making a certain prediction, and to increase clinical
436    confidence in this approach (Lipton, 2017). The lesion percentage is a relevant metric as each RFCA
437    lesion is associated with an arrhythmogenic location of the atrial tissue. The lesions are well defined
438    from simulation of 2D LA models in the current study (and known by a clinician when treating a
439    patient) – but the DL model does not learn the locations of the ablation lesions during training. Hence,
440    the ability of the DL model to utilise these (unseen) lesion locations in its predictions of the RFCA
441    strategy from patient MRI provides foundation for the development of interpretable AI. In the future,
442    such AI approaches can provide a clinician with decision support tools that they understand and trust.

443

## 444    5    Funding

448

449

450

## 451 6    References

452  Adebayo, J., Gilmer, J., Muelly, M., Goodfellow, I., Hardt, M., and Kim, B. (2018). Sanity Checks
453      for Saliency Maps. *Adv Neural Inf Process Syst* 2018-December, 9505–9515. doi:
454      10.48550/arxiv.1810.03292.

455  Alhusseini, M. I., Abuzaid, F., Rogers Albert J and Zaman, J. A. B., Baykaner, T., Clopton, P.,
456      Bailis, P., et al. (2020). Machine Learning to Classify Intracardiac Electrical Patterns During
457      Atrial Fibrillation. *Circ. Arrhythm. Electrophysiol.* 13, e008160.

458  Ancona, M., Ceolini, E., Öztireli, C., and Gross, M. (2017). Towards better understanding of
459      gradient-based attribution methods for Deep Neural Networks. *6th International Conference on
460      Learning Representations, ICLR 2018 - Conference Track Proceedings*. doi:
461      10.48550/arxiv.1711.06104.

462  Bass, C., da Silva, M., Sudre, C., Williams, L. Z. J., Sousa, H. S., Tudosiu, P.-D., et al. (2022).
463      ICAM-reg: Interpretable Classification and Regression with Feature Attribution for Mapping
464      Neurological Phenotypes in Individual Scans. *IEEE Trans Med Imaging*, 1–1. doi:
465      10.1109/TMI.2022.3221890.

466  Bell, A., Solano-Kamaiko, I., Nov, O., and Stoyanovich, J. (2022). It's Just Not That Simple: An
467      Empirical Study of the Accuracy-Explainability Trade-off in Machine Learning for Public
468      Policy. *ACM International Conference Proceeding Series* 22, 248–266. doi:
469      10.1145/3531146.3533090.

470  Chen, S. A., Hsieh, M. H., Tai, C. T., Tsai, C. F., Prakash, V. S., Yu, W. C., et al. (1999). Initiation
471      of atrial fibrillation by ectopic beats originating from the pulmonary veins: electrophysiological
472      characteristics, pharmacological responses, and effects of radiofrequency ablation. *Circulation*
473      100, 1879–1886.

474  Chubb, H., Karim, R., Roujol, S., Nuñez-Garcia, M., Williams, S. E., Whitaker John and Harrison, J.,
475      et al. (2018). The reproducibility of late gadolinium enhancement cardiovascular magnetic
476      resonance imaging of post-ablation atrial scar: a cross-over study. *Journal of Cardiovascular
477      Magnetic Resonance* 20.

478  Chugh, S. S., Havmoeller, R., Narayanan, K., Singh, D., Rienstra, M., Benjamin, E. J., et al. (2014).
479      Worldwide epidemiology of atrial fibrillation: a Global Burden of Disease 2010 Study.
480      *Circulation* 129, 837–847.

481  Fenton, F., and Karma, A. (1998). Vortex dynamics in three-dimensional continuous myocardium
482      with fiber rotation: Filament instability and fibrillation. *Chaos* 8, 20–47.

483  Firouznia, M., Feeny, A. K., Labarbera, M. A., Mchale, M., Cantlay, C., Kalfas, N., et al. (2021).
484      Machine Learning-Derived Fractal Features of Shape and Texture of the Left Atrium and
485      Pulmonary Veins From Cardiac Computed Tomography Scans Are Associated With Risk of
486      Recurrence of Atrial Fibrillation Postablation. *Circ Arrhythm Electrophysiol* 14, 329–336. doi:
487      10.1161/CIRCEP.120.009265.

488  Fukumoto, K., Habibi, M., Gucuk Ipek, E., Khurram, I. M., Zimmerman, S. L., Zipunnikov, V., et al.
489      (2015). Comparison of preexisting and ablation-induced late gadolinium enhancement on left

490    atrial magnetic resonance imaging. *Heart Rhythm* 12, 668–672. doi:
491        10.1016/J.HRTHM.2014.12.021.

492    Graziani, M., Palatnik de Sousa, I., Vellasco, M. M. B. R., Costa da Silva, E., Müller, H., and
493        Andrearczyk, V. (2021). Sharpening Local Interpretable Model-Agnostic Explanations
494        for Histopathology: Improved Understandability and Reliability. *Lecture Notes in Computer*
495        *Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in*
496        *Bioinformatics)* 12903 LNCS, 540–549. doi: 10.1007/978-3-030-87199-4_51/FIGURES/5.

497    Hart, R. G., and Halperin, J. L. (2001). Atrial fibrillation and stroke : concepts and controversies.
498        *Stroke* 32, 803–808.

499    Kheirkhahan, M., Baher, A., Goldooz, M., Kholmovski, E. G., Morris, A. K., Csecs, I., et al. (2020).
500        Left atrial fibrosis progression detected by LGE-MRI after ablation of atrial fibrillation. *Pacing*
501        *Clin Electrophysiol* 43, 402–411. doi: 10.1111/PACE.13866.

502    Kim, J. Y., Kim, Y., Oh, G.-H., Kim, S. H., Choi, Y., Hwang, Y., et al. (2020). A deep learning
503        model to predict recurrence of atrial fibrillation after pulmonary vein isolation. *J. Interv. Card.*
504        *Electrophysiol.* 21, 1–7.

505    Kingma, D. P., and Ba, J. (2014). Adam: A Method for Stochastic Optimization.

506    Kokhlikyan, N., Miglani, V., Martin, M., Wang, E., Alsallakh, B., Reynolds, J., et al. (2020).
507        Captum: A unified and generic model interpretability library for PyTorch.

508    Liesbet Geris, Cécile F. Rousseau, Jerome Noailly, Payman Afshari, Michaël Auffret, Wen-Yang
509        Chu, et al. (n.d.). The Role Of Artificial Intelligence Within In Silico Medicine. Available at:
510        https://www.vph-institute.org/upload/ai-in-health-white-paper_6331c4e3c60cb.pdf [Accessed
511        November 28, 2022].

512    Lipton, Z. C. (2017). The Doctor Just Won't Accept That! doi: 10.48550/arxiv.1711.08037.

513    Liu, C.-M., Chang, S.-L., Chen, H.-H., Chen, W.-S., Lin, Y.-J., Lo, L.-W., et al. (2020). The Clinical
514        Application of the Deep Learning Technique for Predicting Trigger Origins in Patients With
515        Paroxysmal Atrial Fibrillation With Catheter Ablation. *Circ. Arrhythm. Electrophysiol.* 13,
516        e008518.

517    Luongo, G., Azzolin, L., Schuler, S., Rivolta, M. W., Almeida, T. P., Martínez, J. P., et al. (2021).
518        Machine learning enables noninvasive prediction of atrial fibrillation driver location and acute
519        pulmonary vein ablation success using the 12-lead ECG. *Cardiovasc Digit Health J* 2, 126–136.
520        doi: 10.1016/J.CVDHJ.2021.03.002.

521    Magesh, P. R., Myloth, R. D., and Tom, R. J. (2020). An Explainable Machine Learning Model for
522        Early Detection of Parkinson's Disease using LIME on DaTSCAN Imagery. *Comput Biol Med*
523        126, 104041. doi: 10.1016/J.COMPBIOMED.2020.104041.

524    Mahapatra, D., Ge, Z., and Reyes, M. (2022). Self-Supervised Generalized Zero Shot Learning For
525        Medical Image Classification Using Novel Interpretable Saliency Maps. *IEEE Trans Med*
526        *Imaging*. doi: 10.1109/TMI.2022.3163232.

527 Mertes, S., Huber, T., Weitz, K., Heimerl, A., and André, E. (2022). GANterfactual—Counterfactual
528     Explanations for Medical Non-experts Using Generative Adversarial Learning. *Front Artif Intell*
529     5, 825565. doi: 10.3389/FRAI.2022.825565/FULL.

530 Morgan, R., Colman, M. A., Chubb, H., Seemann, G., and Aslanidi, O. v (2016). Slow Conduction in
531     the Border Zones of Patchy Fibrosis Stabilizes the Drivers for Atrial Fibrillation: Insights from
532     Multi-Scale Human Atrial Modeling. *Front. Physiol.* 7, 474.

533 Mourby, M., Ó Cathaoir, K., and Collin, C. B. (2021). Transparency of machine-learning in
534     healthcare: The GDPR &amp; European health law. *Computer Law & Security Review* 43,
535     105611. doi: 10.1016/j.clsr.2021.105611.

536 Muffoletto, M., Fu, X., Roy, A., Varela Marta and Bates, P. A., and Aslanidi, O. v (2019).
537     Development of a Deep Learning Method to Predict Optimal Ablation Patterns for Atrial
538     Fibrillation. in *2019 IEEE Conference on Computational Intelligence in Bioinformatics and*
539     *Computational Biology (CIBCB)*, 1–4.

540 Muffoletto, M., Qureshi, A., Zeidan, A., Muizniece, L., Fu, X., Zhao, J., et al. (2021). Toward
541     Patient-Specific Prediction of Ablation Strategies for Atrial Fibrillation Using Deep Learning.
542     *Front Physiol* 12. doi: 10.3389/FPHYS.2021.674106.

543 Nagel, C., Luongo, G., Azzolin, L., Schuler, S., Dössel, O., and Loewe, A. (2021). Non-Invasive and
544     Quantitative Estimation of Left Atrial Fibrosis Based on P Waves of the 12-Lead ECG—A
545     Large-Scale Computational Study Covering Anatomical Variability. *J Clin Med* 10, 1797. doi:
546     10.3390/JCM10081797.

547 Narayan, S. M., Krummen, D. E., and Rappel, W. J. (2012a). Clinical mapping approach to diagnose
548     electrical rotors and focal impulse sources for human atrial fibrillation. *J Cardiovasc*
549     *Electrophysiol* 23, 447–454. doi: 10.1111/J.1540-8167.2012.02332.X.

550 Narayan, S. M., Krummen, D. E., Shivkumar, K., Clopton, P., Rappel, W. J., and Miller, J. M.
551     (2012b). Treatment of Atrial Fibrillation by the Ablation of Localized Sources: CONFIRM
552     (Conventional Ablation for Atrial Fibrillation With or Without Focal Impulse and Rotor
553     Modulation) Trial. *J Am Coll Cardiol* 60, 628–636. doi: 10.1016/J.JACC.2012.05.022.

554 Oketani, N., Ichiki, H., Iriki, Y., Okui, H., Ryuichi, M., Fuminori, N., et al. (2012). Catheter ablation
555     of atrial fibrillation guided by complex fractionated atrial electrogram mapping with or without
556     pulmonary vein isolation. *J Arrhythm* 28, 311–323.

557 Paszke, A., Gross, S., Massa, F., Lerer Adam and Bradbury, J., Chanan, G., Killeen, T., et al. (2019).
558     PyTorch: An Imperative Style, High-Performance Deep Learning Library. in *Advances in*
559     *Neural Information Processing Systems*, eds. H. Wallach, H. Larochelle, A. Beygelzimer, F. d
560     Alché-Buc, E. Fox, and R. Garnett (Curran Associates, Inc.).

561 Patel, M. I., Singla, S., Ali Mattathodi, R. A., Sharma, S., Gautam, D., and Kundeti, S. R. (2021).
562     Simulating Realistic MRI variations to Improve Deep Learning model and visual explanations
563     using GradCAM. *Proceedings - 2021 IEEE International Conference on Cloud Computing in*
564     *Emerging Markets, CCEM 2021*, 1–8. doi: 10.1109/CCEM53267.2021.00011.

565 Popescu, D. M., Shade, J. K., Lai, C., Aronis, K. N., Ouyang, D., Moorthy, M. V., et al. (2022).
566     Arrhythmic sudden death survival prediction using deep learning analysis of scarring in the
567     heart. *Nature cardiovascular research* 1, 334. doi: 10.1038/S44161-022-00041-9.

568 Qureshi, A., Roy, A., Chubb, H., de Vecchi, A., and Aslanidi, O. (2020). Investigating Strain as a
569     Biomarker for Atrial Fibrosis Quantified by Patient Cine MRI Data. in *2020 Computing in*
570     *Cardiology*, 1–4.

571 Raschka, S. (2018). Model Evaluation, Model Selection, and Algorithm Selection in Machine
572     Learning.

573 Ribeiro, M. T., Singh, S., and Guestrin, C. (2016). Why should I trust you? in *Proceedings of the*
574     *22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining* (New
575     York, NY, USA: ACM).

576 Rodrigo, M., Alhusseini, M. I., Rogers, A. J., Krittanawong, C., Thakur, S., Feng, R., et al. (2022).
577     Atrial fibrillation signatures on intracardiac electrograms identified by deep learning. *Comput*
578     *Biol Med* 145. doi: 10.1016/J.COMPBIOMED.2022.105451.

579 Roney, C. H., Bayer, J. D., Cochet, H., Meo, M., Dubois, R., Jaïs, P., et al. (2018). Variability in
580     pulmonary vein electrophysiology and fibrosis determines arrhythmia susceptibility and
581     dynamics. *PLoS Comput Biol* 14. doi: 10.1371/JOURNAL.PCBI.1006166.

582 Roney, C. H., Beach, M. L., Mehta, A. M., Sim, I., Corrado, C., Bendikas, R., et al. (2020). In silico
583     Comparison of Left Atrial Ablation Techniques That Target the     Anatomical, Structural, and
584     Electrical Substrates of Atrial Fibrillation. *Front. Physiol.* 11, 1145. doi:
585     10.3389/fphys.2020.572874.

586 Roney, C. H., Sim, I., Yu, J., Beach, M., Mehta, A., Alonso Solis-Lemus, J., et al. (2022). Predicting
587     Atrial Fibrillation Recurrence by Combining Population Data and Virtual Cohorts of Patient-
588     Specific Left Atrial Models. *Circ Arrhythm Electrophysiol* 15, e010253. doi:
589     10.1161/CIRCEP.121.010253.

590 Roy, A., Varela, M., and Aslanidi, O. (2018). Image-Based Computational Evaluation of the Effects
591     of Atrial Wall Thickness and Fibrosis on Re-entrant Drivers for Atrial Fibrillation. *Front*
592     *Physiol* 9. doi: 10.3389/FPHYS.2018.01352.

593 Roy, A., Varela, M., Chubb, H., MacLeod Robert and Hancox, J. C., Schaeffter, T., and Aslanidi, O.
594     (2020). Identifying locations of re-entrant drivers from patient-specific distribution of fibrosis in
595     the left atrium. *PLoS Comput. Biol.* 16, e1008086.

596 Salahuddin, Z., Woodruff, H. C., Chatterjee, A., and Lambin, P. (2022). Transparency of deep neural
597     networks for medical image analysis: A review of interpretability methods. *Comput Biol Med*
598     140, 105111. doi: 10.1016/J.COMPBIOMED.2021.105111.

599 Selvaraju, R. R., Cogswell, M., das Abhishek and Vedantam, R., Parikh, D., and Batra, D. (2017).
600     Grad-cam: Visual explanations from deep networks via gradient-based localization. in
601     *Proceedings of the IEEE international conference on computer vision*, 618–626.
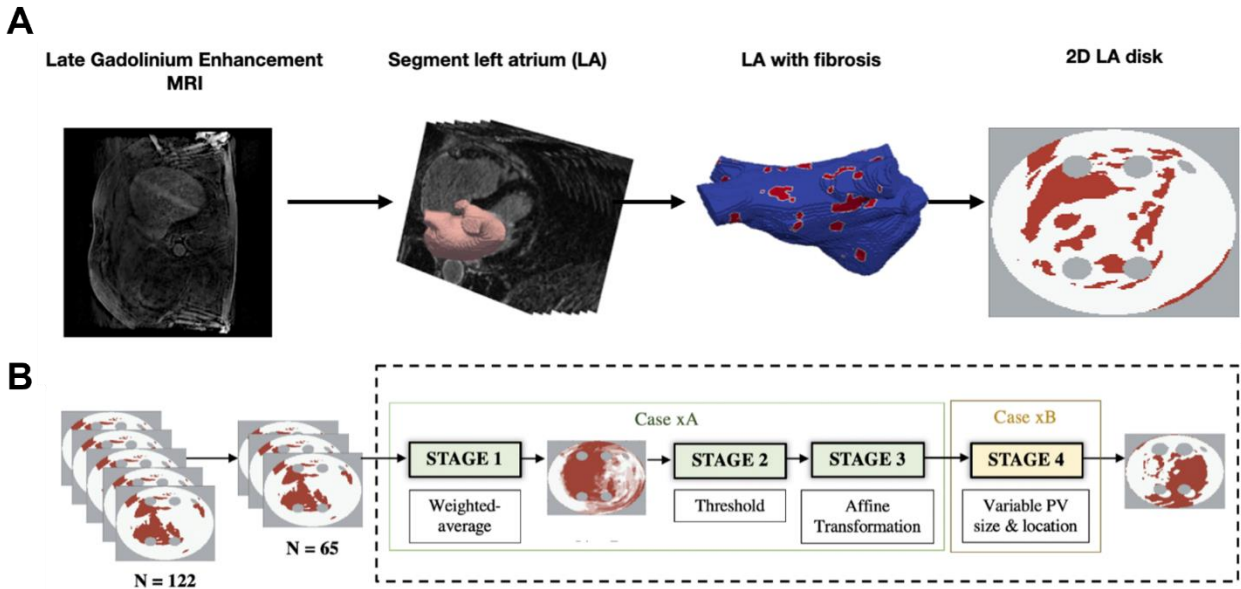
602  Tobón, C., Palacio, L. C., Duque, J. E., Cardona, E. A., Ugarte, J. P., Orozco-Duque, A., et al.
603      (2014). Simple ablation guided by ApEn mapping in a 2D model during permanent atrial
604      fibrillation. in *Computing in Cardiology 2014*, 1029–1032.

605  Townsend, C. M., and Sabiston, D. C. (2001). *Sabiston Review of Surgery*. Saunders.

606  Varela, M., Bisbal, F., Zacur, E., Berruezo, A., Aslanidi, O. v., Mont, L., et al. (2017a). Novel
607      Computational Analysis of Left Atrial Anatomy Improves Prediction of Atrial Fibrillation
608      Recurrence after Ablation. *Front Physiol* 8. doi: 10.3389/FPHYS.2017.00068.

609  Varela, M., Morgan, R., Theron, A., DIllon-Murphy, D., Chubb, H., Whitaker, J., et al. (2017b).
610      Novel MRI Technique Enables Non-Invasive Measurement of Atrial Wall Thickness. *IEEE
611      Trans Med Imaging* 36, 1607–1614. doi: 10.1109/TMI.2017.2671839.

612  Wang, Y., Xu, Y., Ling, Z., Chen, W., Su, L., Du, H., et al. (2017). GW28-e1219 Radiofrequency
613      catheter ablation for paroxysmal atrial fibrillation: over 3-year follow-up outcome. *J Am Coll
614      Cardiol* 70, C126.

615  Watson, D. S., Krutzinna, J., Bruce, I. N., Griffiths, C. E., McInnes, I. B., Barnes, M. R., et al.
616      (2019). Clinical applications of machine learning algorithms: beyond the black box. *BMJ* 364,
617      l886.

618  Williams, S. E., Tobon-Gomez, C., Zuluaga Maria A and Chubb, H., Butakoff, C., Karim, R.,
619      Ahmed, E., et al. (2017). Standardized unfold mapping: a technique to permit left atrial regional
620      data display and analysis. *J. Interv. Card. Electrophysiol.* 50, 125–131.

621  Yubing, W., Yanping, X., Zhiyu, L., Weijie, C., Li, S., Huaan, D., et al. (2018). Long-term outcome
622      of radiofrequency catheter ablation for persistent atrial fibrillation. *Medicine* 97, e11520.

623  Zafar, M. R., and Khan, N. (2021). Deterministic Local Interpretable Model-Agnostic Explanations
624      for Stable Explainability. *Mach Learn Knowl Extr* 3, 525–541.

625  Zeiler, M. D., and Fergus, R. (2014). Visualizing and understanding convolutional networks. *Lecture
626      Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and
627      Lecture Notes in Bioinformatics)* 8689 LNCS, 818–833. doi: 10.1007/978-3-319-10590-1_53.

628  Zhao, J., Aslanidi, O., Kuklik, P., Lee, G., Tse, G., Niederer, S., et al. (2019). Editorial: Recent
629      Advances in Understanding the Basic Mechanisms of Atrial    Fibrillation Using Novel
630      Computational Approaches. *Front. Physiol.* 10, 1065. doi: 10.3389/fphys.2019.01065.

631  Zolotarev, A. M., Hansen, B. J., Ivanova, E. A., Helfrich, K. M., Li, N., Janssen, P. M. L., et al.
632      (2020). Optical Mapping-Validated Machine Learning Improves Atrial Fibrillation Driver
633      Detection by Multi-Electrode Mapping. *Circ Arrhythm Electrophysiol* 13, E008249. doi:
634      10.1161/CIRCEP.119.008249.
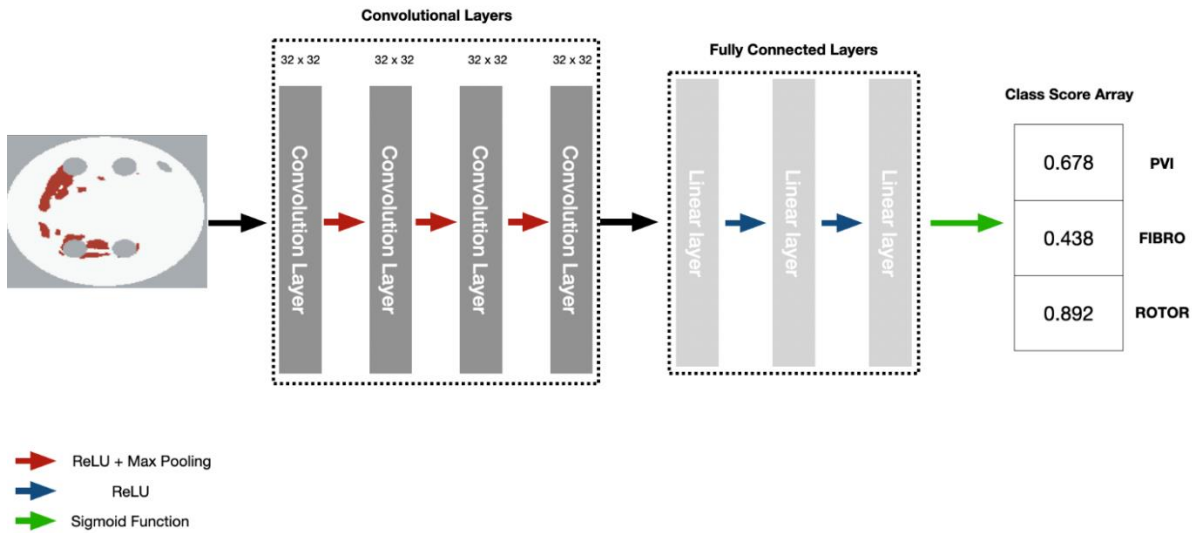
635

636

637    **7    Figures and Tables**

638



639

**Figure 1. Diagram of MRI-derived 2D LA tissue disk. A.** Workflow of 2D LA tissue generation
pipeline. The figure illustrates the process of how the 2D LA tissue models are obtained from LGE
MRI by LA segmentation, thresholding fibrosis from healthy tissue and mapping onto 2D LA tissues.
**B.** Workflow for generating synthetic tissues. 65 tissues were randomly selected from the total dataset
of 122 real tissues. These 65 tissues were used to generate the synthetic images by iterating overstages
1 to 4 (199 times) to create a virtual cohort of 199 tissues. 'Case xA' denotes the combination of data
augmentation techniques used to create the synthetic fibrosis distributions. 'Case xB' determines how
the PV sizes and locations were varied from those in the standardised discs.

648

649

650

**Figure 2.** Diagram of CNN with parameters to predict RFCA simulation strategy success from 2D LA tissue.

| Ablation Strategy | AUC | Recall | Precision | F1-Score |
|---|---|---|---|---|
| **PVI** | $0.78 \pm 0.03$ | $0.35 \pm 0.07$ | $0.68 \pm 0.28$ | $0.42 \pm 0.06$ |
| **FIBRO** | $0.92 \pm 0.02$ | $0.89 \pm 0.03$ | $0.82 \pm 0.02$ | $0.85 \pm 0.01$ |
| **ROTOR** | $0.77 \pm 0.02$ | $0.93 \pm 0.04$ | $0.76 \pm 0.02$ | $0.84 \pm 0.01$ |

**Table 1.** Mean area under the receiver operating characteristic curve (AUC) score, recall, precision and F1-score on independent hold-out test set (with standard deviation) for each RFCA strategy.

| Ablation Strategy | MRI Derived Data | MRI Derived + Synthetic Data |
|---|---|---|
| **PVI** | $0.67 \pm 0.03$ | $0.78 \pm 0.04$ |
| **FIBRO** | $0.85 \pm 0.02$ | $0.92 \pm 0.02$ |
| **ROTOR** | $0.62 \pm 0.05$ | $0.77 \pm 0.02$ |

**Table 2.** Mean AUC score on independent hold-out test set (with standard deviation) for each RFCA strategy and type of data

660

| Ablation Strategy | MRI Derived Data | | | | MRI Derived + Synthetic Data | | | |
|---|---|---|---|---|---|---|---|---|
| | **AUC** | **Recall** | **Precision** | **F1 Score** | **AUC** | **Recall** | **Precision** | **F1 Score** |
| **PVI** | 0.67 ± 0.03 | 0 | 1.0 | 0 | 0.78 ± 0.03 | 0.35 ± 0.07 | 0.68 ± 0.28 | 0.42 ± 0.06 |
| **FIBRO** | 0.85 ± 0.02 | 0.75 ± 0.08 | 0.70 ± 0.03 | 0.72 ± 0.04 | 0.92 ± 0.02 | 0.89 ± 0.03 | 0.82 ± 0.02 | 0.85 ± 0.01 |
| **ROTOR** | 0.62 ± 0.05 | 0.99 ± 0.02 | 0.64 ± 0.01 | 0.78 ± 0.02 | 0.77 ± 0.02 | 0.93 ± 0.04 | 0.76 ± 0.02 | 0.84 ± 0.01 |

661 **Table 3.** Mean AUC, recall, precision and F1 score (with standard deviation) of DL model trained with
662 real data only and with synthetic and real data from a leave-one-out cross-validation on a hold-out test
663 (~20% of the respective dataset).

664

665

666

667

668

669

670

671

672

673

674

675

| Ablation Strategy | Method | Lesion Percentage | IoU | NAT Percentage |
|---|---|---|---|---|
| **PVI** | LIME | 0.44 ± 0.24 | **0.077 ± 0.023** | **0.32 ± 0.24** |
| | Occlusions | **0.55 ± 0.15** | 0.065 ± 0.17 | 0.57 ± 0.15 |
| | GradCAM | 0.47 ± 0.17 | 0.063 ± 0.029 | 0.60 ± 0.12 |
| **FIBRO** | LIME | 0.57 ± 0.19 | 0.18 ± 0.09 | 0.47 ± 0.27 |
| | Occlusions | 0.45 ± 0.14 | 0.19 ± 0.11 | 0.38 ± 0.20 |
| | GradCAM | **0.62 ± 0.25** | **0.26 ± 0.11** | **0.27 ± 0.16** |
| **ROTOR** | LIME | 0.62 ± 0.16 | 0.12 ± 0.07 | 0.63 ± 0.25 |
| | Occlusions | 0.53 ± 0.16 | 0.14 ± 0.06 | 0.36 ± 0.16 |
| | GradCAM | **0.71 ± 0.13** | **0.20 ± 0.08** | **0.25 ± 0.06** |

676 **Table 4.** Mean lesion percentage, NAT percentage, IoU of the informative region and ablation lesions
677 with errors (standard deviation) for each FA map method and RFCA strategy.

678

679

680

**Figure 3.** Diagram of 2D LA tissues with highlighted feature attribution maps. White areas in the 2D tissues are healthy tissue and red areas are fibrosis. Ablation lesion locations known from simulations are shown (yellow) for all three RFCA strategies, along with respective FA maps for LIME, GradCAM and occlusions and highlighted thresholded informative regions (translucent green). Same colour scheme in used in Figures 7 and 8 below.

**Figure 4.** Boxplot of Jacquard index (IoU) for each FA method (GradCAM, LIME and Occlusions) and RFCA strategy (PVI, FIBRO and ROTOR) on the hold-out test set.



**Figure 5.** Boxplot of lesion percentage for each FA method (GradCAM, LIME and Occlusions) and RFCA strategy (PVI, FIBRO and ROTOR) on the hold-out test set.

708



709

**Figure 6.** Boxplot of NAT percentage for each FA method (GradCAM, LIME and Occlusions) and RFCA strategy (PVI, FIBRO and ROTOR) on the hold-out test set.

712

**Figure 7.** Correct and incorrect classification examples of FA maps (LIME, GradCAM and occlusions) for FIBRO.

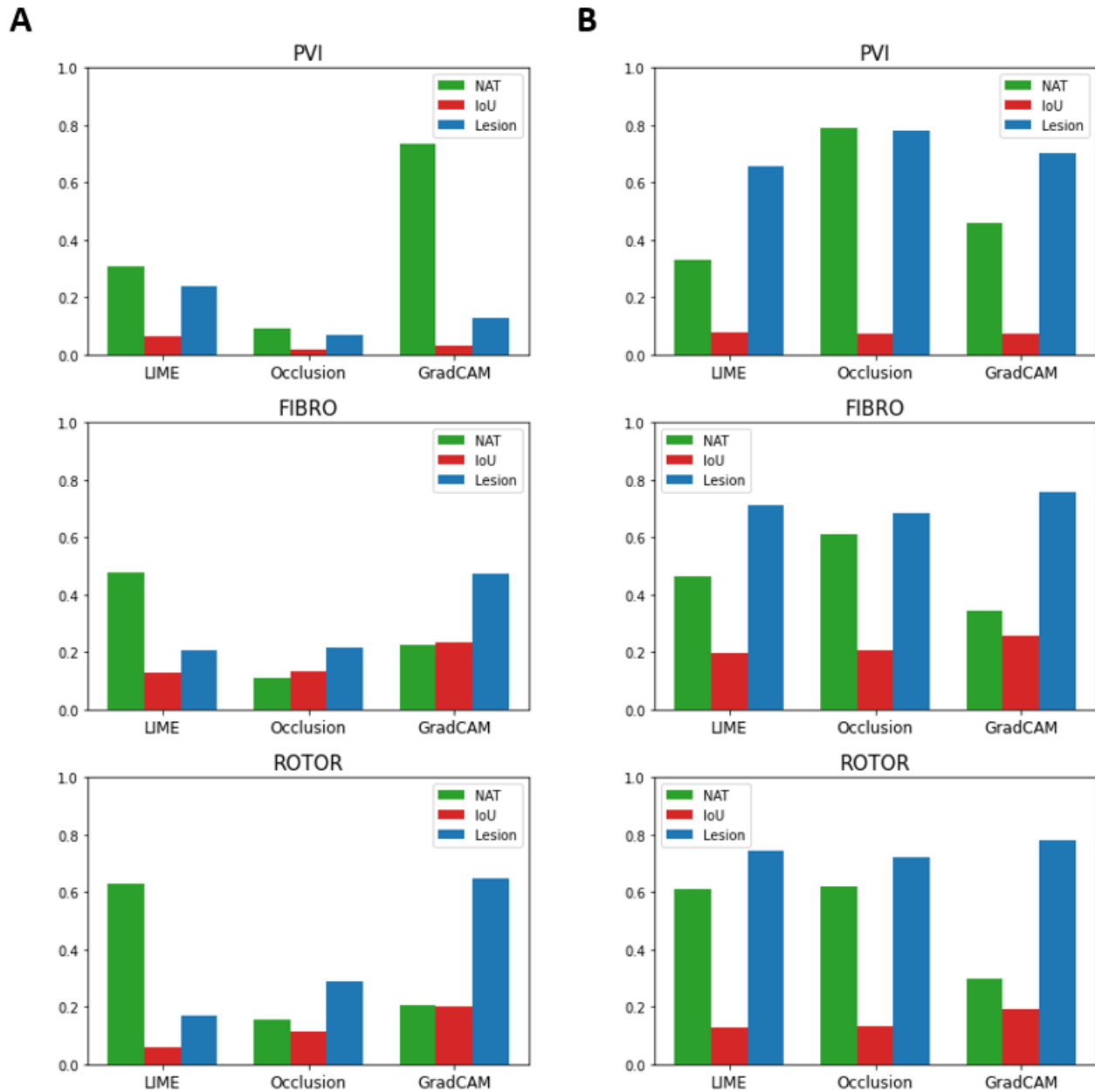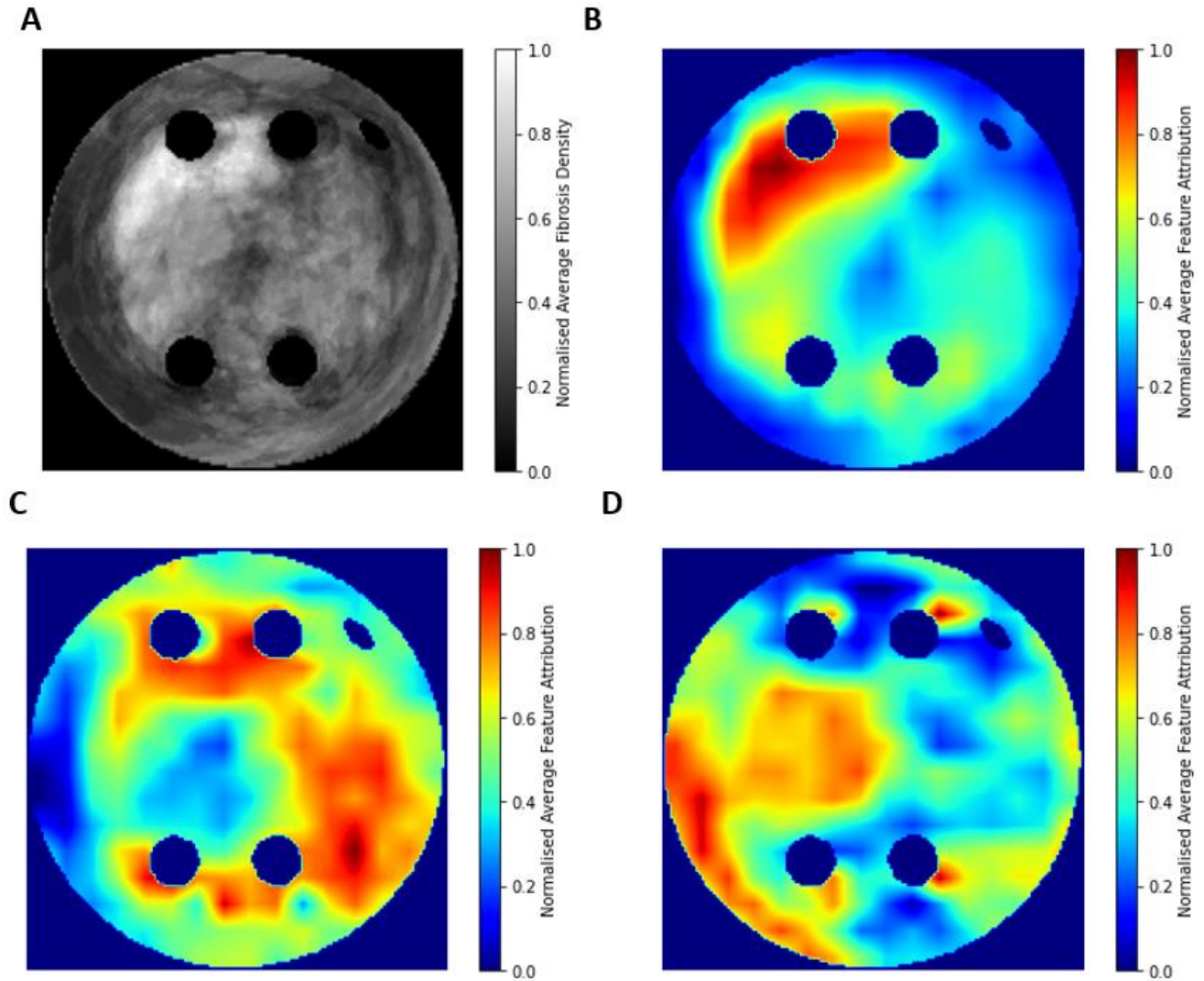**Figure 8.** Correct and incorrect classification examples of FA maps (LIME, GradCAM and occlusions) for ROTOR.

**Figure 9.** IoU, lesion and NAT percentage values for each interpretability method and ablation strategy with altered informative region threshold value. **A**. Informative region threshold value 25% above the average FA. **B**. Informative region threshold value 25% below the average FA.
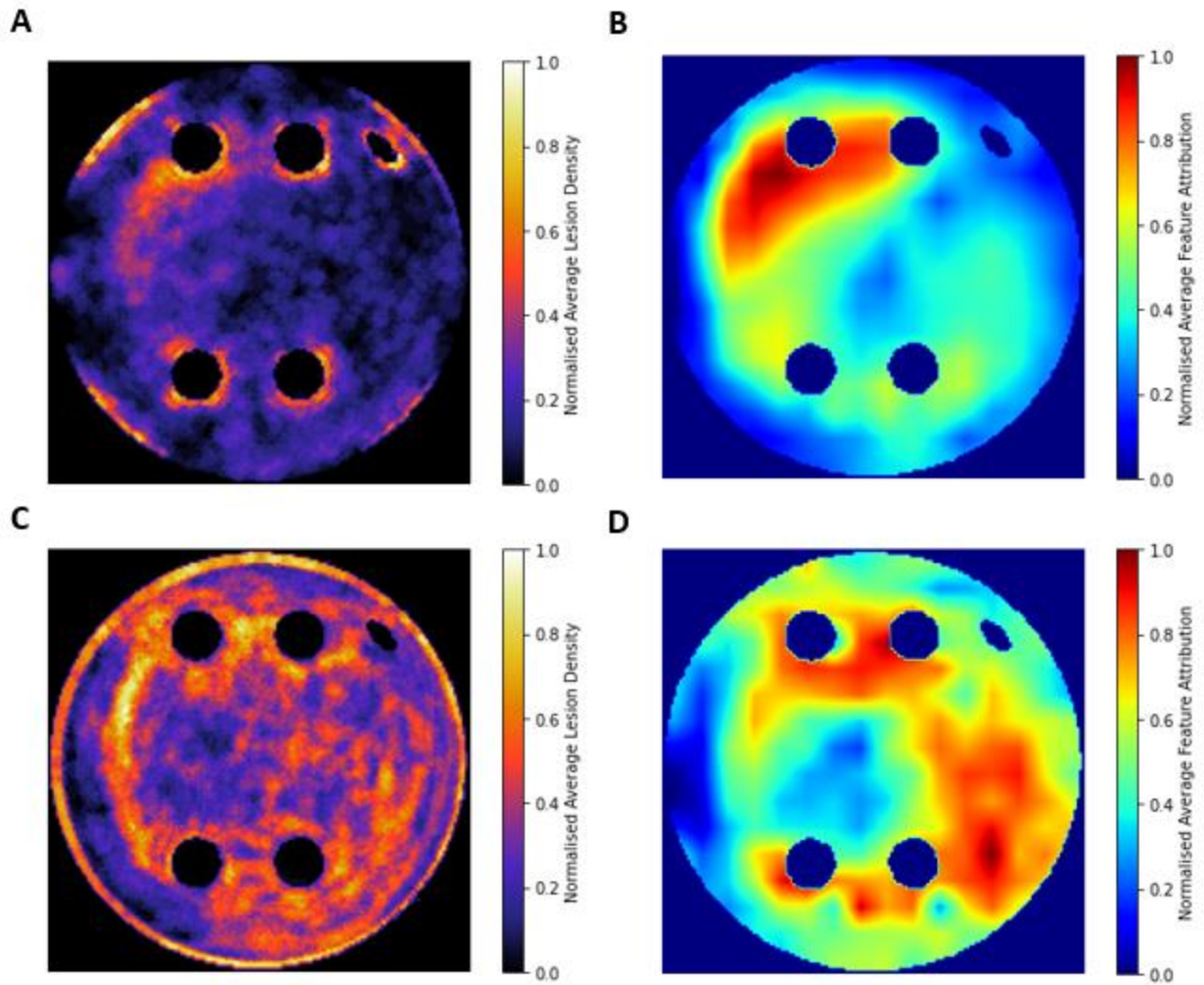
**Figure 10.** Averaged LGE MRI intensities and FA maps on the hold-out test set**. A.** Averaged and normalised LGE MRI intensity in the LA tissue disks. **B**. Averaged and normalised GradCAM FA map for the ROTOR ablation strategy. **C.** Averaged and normalised GradCAM FA map for the FIBRO ablation strategy. **D**. Averaged and normalised GradCAM FA map for the PVI ablation strategy.

**Figure 11.** Averaged and normalised ablation lesions and GradCAM FA maps for FIBRO and ROTOR on the hold-out test set. **A**. Ablation lesions for ROTOR. **B**. FA map for ROTOR. **C**. Ablation lesions for FIBRO. **D**. FA map for FIBRO.