

1 **Diagnostic signature for Heart Failure with Preserved Ejection Fraction (HFpEF): A**  
2 **Machine Learning Approach Using Multi-Modality Electronic Health Record Data**

3 **Short Title:** Diagnosis of HFpEF: A Machine Learning Approach

4 Nazli Farajidavar,<sup>1\*</sup> Kevin O’Gallagher,<sup>1,4\*</sup> Daniel Bean,<sup>1,2,3</sup> Adam Nabeebaccus,<sup>1,4</sup> Rosita

5 Zakeri,<sup>1,4</sup> Daniel Bromage,<sup>1,4</sup> Zeljko Kraljevic,<sup>2</sup> James TH Teo,<sup>4</sup> Richard J Dobson,<sup>1,2,3,5</sup> Ajay M

6 Shah.<sup>1,4</sup>

7 \*joint authors

8

9 <sup>1</sup>King's College London British Heart Foundation Centre of Excellence, School of

10 Cardiovascular Medicine & Sciences, London, UK;

11 <sup>2</sup>Department of Biostatistics and Health Informatics, Institute of Psychiatry, Psychology and

12 Neuroscience, King’s College London, London, UK;

13 <sup>3</sup>Health Data Research UK London, Institute of Health Informatics, University College London,

14 London, U.K

15 <sup>4</sup>King’s College Hospital NHS Foundation Trust, London, UK;

16 <sup>5</sup>NIHR Biomedical Research Centre at South London and Maudsley NHS Foundation Trust and

17 King’s College London, London, UK.

18

19 **Correspondence:** Prof Ajay M Shah, School of Cardiovascular Medicine & Sciences, King’s

20 College London, James Black Centre, 125 Coldharbour Lane, London SE5 9NU, UK. Tel: 0044-

21 207848-5189. Email: [ajay.shah@kcl.ac.uk](mailto:ajay.shah@kcl.ac.uk)

22 **Word count: 3414**

23 **ABSTRACT**

24 **Aims:** Heart failure with preserved ejection fraction (HFpEF) is thought to be highly prevalent  
25 yet remains underdiagnosed. We sought to develop a data-driven  
26 diagnostic model to predict from electronic health records (EHR) the likelihood  
27 of HFpEF among patients with unexplained dyspnea and preserved left ventricular EF.

28 **Methods & Results:** The derivation cohort comprised patients with dyspnea and  
29 echocardiography results. Structured and unstructured data were extracted using an automated  
30 informatics pipeline. Patients were retrospectively diagnosed as HFpEF (cases), non-HF  
31 (control cohort I), or HF with reduced EF (HFrEF; control cohort II). The ability of clinical  
32 parameters and investigations to discriminate cases from controls was evaluated by extreme  
33 gradient boosting. A likelihood scoring system was developed and validated in a separate test  
34 cohort.

35 The derivation cohort included 1585 consecutive patients: 133 cases of HFpEF (9%), 194 non-  
36 HF cases (Control cohort I) and 1258 HFrEF cases (Control cohort II). Two HFpEF diagnostic  
37 signatures were derived, comprising symptoms, diagnoses and investigation results. A final  
38 prediction model was generated based on the averaged likelihood scores from these two models.  
39 In a validation cohort consisting of 269 consecutive patients (with 66 HFpEF cases (24.5%)), the  
40 diagnostic power of detecting HFpEF had an AUROC of 90% ( $P < 0.001$ ) and average precision  
41 (AP) of 74%.

42 **Conclusion:** This diagnostic signature enables discrimination of HFpEF from non-cardiac  
43 dyspnea or HFrEF from EHR and can assist in the diagnostic evaluation in patients  
44 with unexplained dyspnea.

45

46 **Key words:** HFpEF, machine learning, dyspnea

## 47 INTRODUCTION

48 Heart Failure with preserved ejection fraction (HFpEF) is a highly prevalent yet under-diagnosed  
49 clinical syndrome[1, 2]. The hallmarks are the signs and symptoms of heart failure (HF) and a  
50 preserved left ventricular ejection fraction (LVEF). HFpEF is thought to be underpinned by  
51 structural and functional abnormalities of both the heart and vasculature. Patients with HFpEF  
52 typically display diastolic dysfunction[3, 4] and other abnormalities such as vascular  
53 stiffening[5] and impaired ventricular-vascular coupling[6-10]. Unlike HF with reduced Ejection  
54 Fraction (HFrEF), no evidence-based therapies are available for HFpEF[11-13]. This may in  
55 part reflect the heterogeneity of HFpEF pathophysiology as well as issues of clinical trial  
56 design[13-15].

57 While the diagnosis of HFpEF is straightforward in acutely decompensated patients, stable  
58 euvolemic patients present a greater challenge[16]. Exertional dyspnea and fatigue are non-  
59 specific symptoms that occur in many other conditions, including obesity and physical  
60 deconditioning. Expert transthoracic echocardiography (ideally with exercise) or invasive cardiac  
61 catheterization to document raised LV filling pressures may not be immediately available to the  
62 non-specialist. A recent study found that among more than 44,000 community-based patients  
63 likely to have HF, only 50% had a documented LVEF[17]. Furthermore, those eventually  
64 diagnosed as having HFpEF required many more pre-diagnosis investigations and consultations  
65 than HFrEF patients.

66 In previous epidemiological studies, identification and extraction of HFpEF cases from  
67 Electronic Health Records (EHR) has typically relied on diagnostic codes, additional medical  
68 record abstraction, and/or adjudication based on various expert criteria e.g. European Society of

69 Cardiology criteria[18]. The EHR is however increasingly amenable to rapid and automated  
70 extraction of multiple clinical parameters, including the use of advanced natural language  
71 processing (NLP) algorithms to identify clinical concepts recorded in the unstructured text[19-  
72 21].

73 The aim of this study was to extract and analyze data from the EHR to develop an automated  
74 approach to identify patients likely to have HFpEF.

75

## 76 **METHODS**

### 77 *Approvals*

78 This project was conducted under London South East Research Ethics Committee approval  
79 (reference 18/LO/2048) granted to the King's Electronic Records Research Interface (KERRI),  
80 project ID 202020201.

### 81 *Derivation Cohort*

82 We performed a retrospective study using de-identified data of patients attending King's College  
83 Hospital NHS Foundation Trust (KCH) in London (UK) between 2000 and 2019. We focused on  
84 patients who had undergone echocardiography as part of their inpatient or outpatient evaluation.  
85 With this starting point, a number of different patient cohorts were derived based on the LVEF,  
86 confirmed or possible HF, symptoms of dyspnea, and NT-proBNP (or BNP) level (see  
87 **Supplementary materials Sections I and II**). We identified confirmed HFpEF cases and two  
88 control cohorts: those with no evidence of HF (non-HF, Control cohort I) and those with HFrEF  
89 (Control cohort II). HFpEF cases were defined as patients with a preserved LVEF  $\geq 50\%$  (with  
90 no evidence of LVEF  $< 50\%$  at any stage), a confirmed diagnosis of HF based on discharge  
91 ICD10 codes I50.0, I50.1 or I50.9, dyspnea, and a raised NT-proBNP or BNP level (according to  
92 age-specific thresholds), in accordance with ESC diagnostic criteria[18]. Patients with valvular  
93 heart disease (ICD10 codes I05-I09 and I35) were excluded.

### 94 *Test Cohorts*

95 We generated 4 test cohorts from patients who lacked at least one of the above diagnostic  
96 features for a confirmed diagnosis of HFpEF (see **Supplementary Table S1 and Flowchart S1**).  
97 We randomly sampled 100 patients from each of these four test subsets for analysis and removed

98 samples where the clinical annotations disagreed or there was more than 70% missingness in  
99 signature predictors, leaving 269 in total.

### 100 *Data extraction and evaluation*

101 Clinical and demographic data were retrieved from the structured and unstructured components  
102 of the EHR using the CogStack informatics platform[20]. Automated parsing of the EHR was  
103 achieved with a state-of-the-art enterprise search and well-validated natural language processing  
104 (NLP) tools, including MedCAT[22] and the Unified Medical Language System repository[23]  
105 as previously used by our group.[24] Clinical term extraction was restricted to concepts which  
106 represent clinical findings, diseases (apart from HF), medications, and signs and symptoms. This  
107 was linked to searches of structured data from an internal database containing echocardiographic  
108 data and ICD codes. Continuous variables were cleaned prior to cohort selection; e.g. conversion  
109 of text references of LVEF to numerical values and removal of measurement outliers (see  
110 **Supplementary material Section III**). We used both platforms to arbitrate discrepancies in our  
111 derivation dataset as neither source proved to be comprehensive, in line with previous work[20,  
112 21].

113 Echocardiographic data were based on studies performed according to British Society of  
114 Echocardiography guidelines[25] (which are consistent with American and European  
115 guidelines)[26]. Structured data recorded in echocardiography results were boosted with  
116 numerical data reported in the EHR text. Additionally, when appropriate (e.g. patient had  
117 echocardiography but a numerical value for LVEF was not documented) we used a deep learning  
118 model to infer whether the LVEF was preserved based on the echocardiography report (see  
119 **Supplementary materials Section III**).

120 BNP or NT-proBNP results were obtained from samples drawn at any time in the study period  
121 and the maximum value for each subject was used.

122 All cases in the derivation dataset that were identified by the data pipeline as HFpEF were  
123 validated by manual review of the EHR by a cardiologist.

124

### 125 *Potential modeling predictors*

126 A binary diagnostic outcome indicating the presence or absence of HFpEF was considered for  
127 modeling. Potential predictors to be included in a diagnostic signature included those used in  
128 previous HFpEF epidemiological studies[14, 15]. In addition, we adopted a comprehensive  
129 approach that included physiological variables, laboratory results, echocardiographic data and  
130 clinical concept references[27]. Structured data were collected within a two-month temporal  
131 window around the last echocardiography result (or NTproBNP/BNP test result if available).  
132 Unstructured data were analyzed from the entire EHR prior to the date of the echocardiography  
133 result for each patient.

134 We made a second level predictor grouping according to whether the variables were initially  
135 recorded as (a) structured data: demographic and physiological parameters, and laboratory and  
136 echocardiography measurements; or (b) unstructured text in the EHR, extracted via the NLP  
137 platform. We adopted the bag-of-words[28] approach to transform clinical concept annotation  
138 into word vectors for modeling purposes. Concepts which were mentioned in <10% of the  
139 derivation cohort were excluded. Data from the other predictor categories were collected and  
140 imputed prior to training, using the k-nearest neighbor (Scikit-learn python package v0.22) after



141 min-max normalization. Following imputation, data items were rescaled into their original range  
142 to preserve the explainability of the final model.

### 143 *Data modelling, feature selection and validation*

144 We used the tree-based multivariable extreme gradient boosting[29] algorithm (XGBoost, python  
145 package v0.9) for modeling, enabling inclusion of mixed data types and smooth handling of  
146 missing values and sparsity issues. As such, when a value is missing in the sparse predictor  
147 vector, the instance is classified into a default direction (see[29] for further details) that is learnt  
148 as optimal using derivation data.

149 SHAP[30] analysis (SHapley Additive exPlanations; SHAP python package v0.33) was used to  
150 order the predictors according to their prominence in discriminating cases from controls. Once  
151 the full model was created, we took a stepwise forward insertion scheme to include the  
152 more significant variables one at a time, in order to determine the minimal number of predictors  
153 that gave an acceptable performance relative to the use of all predictors. The final predictive  
154 models were trained and evaluated using the obtained optimal subset of predictors.

155 Model validation was undertaken in the test cohorts described earlier, using clinical assessment  
156 criteria from the H<sub>2</sub>FPEF score[16] as a comparator. A random sample of 400 patients from the  
157 test datasets was manually reviewed by two teams each comprising two cardiologists, in order to  
158 validate diagnoses. Any cases of clinician disagreement were removed from the evaluation,  
159 leaving a total of 269 patients in the test datasets (see Results, **Table 1**).

### 160 **Statistical analysis of predictors**

161 Data are presented as mean and standard deviation (SD) or median and interquartile range (IQR)  
162 as appropriate. Differences between cases and controls were evaluated by the Mann-Whitney U  
163 test or unpaired t test, as appropriate. The area under the receiver-operating characteristic curve  
164 (AUROC), F1-score (macro and weighted<sup>2</sup>) and average precision (AP) were used as  
165 performance metrics.

166 A stratified 5-fold cross-validation scheme (to ensure each fold is a good representative of the  
167 whole data in terms of class prevalence) was utilized for feature selection and derivation set  
168 validation. As such, the derivation data was divided into five subsets, four of which were used  
169 for training the model and the final one for validation/testing. The derivation and test subsets  
170 were shuffled until all five subsets were evaluated. The final performance was then reported as  
171 mean and standard deviation of all five tests.

172 The AUROC and AP were used as performance metrics and the Kappa statistic was used to  
173 measure the inter-rater agreement of proposed models. All tests were 2-sided, with  $P < 0.05$   
174 considered significant.

175 To evaluate the generalizability of the model to a new sample, Harrell optimism was calculated  
176 with 1000 boot-strap replicates[31]. To evaluate discrimination power of the proposed model  
177 beyond existing criteria, we compared the model's AUROCs and AP performance against the  
178 recently proposed H<sub>2</sub>FPEF scoring system[16] using the Random Forest (predecessor  
179 to XGBoost).

180 Statistical analyses were performed in Python 3 using SciPy and Scikit-learn packages (v0.22).

181 **Data availability**

182 The data included in the study will not be made available to other researchers due to hospital  
183 information governance regulations. However, we will share our models and the analytical  
184 methods to facilitate the replication of the study on data collected from other hospitals.

185

## 186 RESULTS

187 1854 patients were included in the study of whom 1585 were in the derivation cohort (**Table 1**).  
188 HFpEF patients in the derivation cohort (n=133) were older than those with non-HF or HFrEF,  
189 with a higher proportion of females and a higher BMI. They also had a higher prevalence of  
190 hypertension, atrial fibrillation, diabetes and chronic kidney disease. Systolic and diastolic  
191 pressures were higher in the HFpEF group compared to HFrEF. Patients with HFpEF had lower  
192 end-diastolic and end-systolic volumes and higher septal E/e' ratios than the non-HF control  
193 group.

194

### 195 *Structured, unstructured and combined signatures for HFpEF diagnosis*

196 We initially divided the predictors into two sets based on the source of data being structured data  
197 or clinical concepts and conditions extracted from the unstructured historical EHR (*see*  
198 *Methods*). We excluded the BNP/NT-proBNP assessment data and HF concept references from  
199 both predictor sets to avoid biasing models by information on outcome. Separate  
200 XGBoost models were trained on each predictor set. SHAP analysis was adopted to select the  
201 optimal number of features from each predictor set using five-fold cross-validation. We then  
202 compared the discriminant power of these signatures to distinguish HFpEF cases either from  
203 non-HF patients (Control set I) or HFrEF patients (Control set II).

204 The minimum number of variables required to maintain an acceptable level of performance for  
205 each model were selected (**Figure 1**). Following an early-fusion modeling strategy, we merged  
206 the selected predictors from the two sets of structured and unstructured variables and trained

207 an XGBoost model for discrimination and termed the derived signature as the combined  
208 signature.

209 SHAP analysis to assess feature importance showed that individual predictors had different value  
210 in discrimination of HFpEF versus non-HF or HFrEF (**Figure 2**). For example, dyspnea and  
211 pharmacologic substance were the most prominent predictors in discrimination against non-HF  
212 whereas EF was most important for discrimination against HFrEF. However, many of the  
213 features (e.g. age, patient address) were common to the two groups. The text references to  
214 “patient address” and “pharmacologic substance” were surrogate predictors of the number of  
215 complete hospital admissions. (**Figure 2**).

216 The combined signature model for discrimination of HFpEF from HFrEF showed an enhanced  
217 AUROC performance and F1-measure score as compared to the single-view models in the 5-fold  
218 cross-validation evaluation in our developmental dataset (**Table 2**). The performance  
219 enhancement of the combined model in discriminating HFpEF from non-HF was less significant.  
220 This was due to dominance of the unstructured predictors in this combined signature (see **Figure**  
221 **2** and **Table 3**).

### 222 *Selection of the final model and evaluation in test cohorts*

223 The final model that was used for test evaluations aggregates the HFpEF vs HFrEF and HFpEF  
224 vs non-HF signature likelihood predictions, through an averaging operation. We used this  
225 aggregate model to make predictions on the test sets. **Figure S5** summarises the entire  
226 processing and model training pipeline.

227 To address the distributional variation between training and test cohorts which was caused by  
228 sample selection bias, we used 30% of the test samples (test1: 19, test2: 21, test3: 17, test4: 23)  
229 to retune the models, following the domain adaptation transfer learning technique[32]. Details of  
230 the 30% choice of adaptation set size is included in the **Supplementary materials, Figure S6**.

231 The performance of both proposed base models and the final aggregated model remained robust  
232 in the test cohort as compared to expert clinical consensus, with an AUROC performance of 0.86  
233 (95% CI,  $\pm 0.002$ ) and 0.85 (95% CI,  $\pm 0.001$ ) in HFpEF vs non-HF and HFpEF vs HFrEF  
234 models, respectively and an enhanced aggregate performance of 0.90 (95% CI,  $\pm 0.002$ ) in our  
235 final aggregate model (**Figure 3**).

236 Lastly, we compared the final aggregate model as well as the baseline combined signature  
237 models (discriminating against non-HF or HFrEF) with the recently described H<sub>2</sub>FPEF  
238 model[16]. The AUROC and average precision of both the aggregate model and the individual  
239 baseline models was higher than the H<sub>2</sub>FPEF model (**Table 4**). We additionally used the Cohen's  
240 kappa score to report on the agreement between the predictions made by our proposed models to  
241 better highlight the efficiency of the aggregate model over the individual base models  
242 discriminating HFpEF from non-HF and HFrEF. The positive kappa score of 0.3 indicates a  
243 weak agreement between the two base models. This was expected as the test cohort had lower  
244 availability of clinical assessments compared to the derivation cohort.

245

## 246 **DISCUSSION**

247 In this study, we have developed an automated pipeline for EHR-based data collection,  
248 processing and modeling to identify patients with a high likelihood of HFpEF. We incorporated  
249 multi-modality data, including both structured and unstructured predictors, to generate a disease  
250 diagnostic signature. The proposed signature was validated in a separate cohort of patients and  
251 performed favourably as compared either to expert clinical consensus or the recently proposed  
252 H<sub>2</sub>FPEF score[16].

253 Analysis of the signatures that distinguished HFpEF from non-cardiac causes of dyspnea (non-  
254 HF) revealed anticipated predictors such as atrial fibrillation, hypertension, diabetes mellitus,  
255 kidney failure and obesity, in accordance with previous literature[16]. In addition, surrogate  
256 measures of multiple previous clinical encounters detected by the NLP algorithm as frequent text  
257 references to terms such as “pharmacologic substance” or “patient address” were very useful.  
258 This may reflect the fact that patients with HFpEF may require multiple clinical visits and  
259 investigations, often with different specialities, before a diagnosis is established[17]. Apart from  
260 LVEF itself, features that distinguished HFpEF from HFrEF included age, peripheral edema, and  
261 other echocardiographic measures. An advantage of the approach that we employed may be that  
262 it is unbiased and comprehensive and identifies variables for inclusion in the diagnostic signature  
263 based purely on the results of the objective feature selection process. This may be one reason  
264 why our algorithm outperforms the H<sub>2</sub>FPEF score, which is based on the evaluation of selected  
265 variables rather than a comprehensive unbiased analysis. In this regard, it is of interest that  
266 echocardiographic predictors that contributed to the differentiation of HFpEF from HFrEF  
267 included maximum flow velocity across the aortic valve, aortic insufficiency and LA volume  
268 whereas E/e’ (which is part of the H<sub>2</sub>FPEF score) did not feature in the top 30 predictors.

269 A major underlying problem in efforts to develop or test new treatments for HFpEF is the  
270 difficulty in consistently diagnosing the syndrome[17]. Many different approaches are used in  
271 the literature based on varying criteria published by national and international societies, and  
272 diverse inclusion criteria have been used in clinical trials[33-35]. The problem is compounded by  
273 the likelihood that HFpEF is a heterogenous syndrome in which sub-populations may have  
274 differing underlying pathophysiology and outcomes[14, 15, 33]. The approach we present  
275 enables rapid identification of likely HFpEF cases among which further specific phenotyping  
276 could be performed to refine the diagnosis and potentially test or target defined interventions, or  
277 to identify potential subjects for research studies. Importantly, this approach aims to identify  
278 both compensated and decompensated HFpEF cases, using an automated and data-driven  
279 approach that is effective even where structured data (e.g. NT-proBNP measurements) are  
280 scarce. The approach may be considered complementary to scores such as H<sub>2</sub>FPEF. Our  
281 signature is ideally suited to rapidly identify a large number of possible HFpEF cases from EHR  
282 whereas H<sub>2</sub>FPEF is better suited for use by the clinician evaluating an individual patient who is  
283 suspected to have HFpEF.

284 This study is the first to use SHAP analysis for feature selection in this context. We  
285 comprehensively validated all variations of the derived models in multiple datasets with  
286 underlying variational distributions. We demonstrated a significant improvement in HFpEF  
287 diagnostic performance when discriminating the patients with HFpEF from those with HFrEF or  
288 no HF history. A key strength of our approach is that modeling numerical assessment data  
289 (structured results signature) and EHR concept references separately makes the models  
290 applicable in scenarios where one of these sources of data may be scarce. Moreover, the dual



291 modeling of HFpEF separation from non-HF and HFrEF subjects increases the utility of the  
292 proposed pipeline in distinguishing among a wider group of clinical conditions.

293

### 294 *Limitations*

295 The UMLS clinical concept encoding that was used to extract unstructured observations does not  
296 support distinct encoding of different disease stages and could therefore cause some inaccuracy.

297 In a more general aspect, the *a priori* assumptions that we made to identify definite HFpEF cases

298 in the derivation dataset influenced the characterisation of the cohort. For example, we utilised

299 ICD-10 diagnostic codes in the identification of patients with heart failure. Previous studies

300 have demonstrated inaccuracy in identifying incident heart failure using ICD-10 coding as the

301 sole source[36]. It is possible that such inaccuracy is present in our coding system; however the

302 use of additional features (symptoms, LVEF, BNP/NTproBNP) in case classification mitigates

303 this risk in our study. The inclusion of a raised BNP criterion restricts the cohort to a subgroup

304 of HFpEF subjects, which was evident in test cohorts where many of the subjects did not have

305 BNP measurements. This issue could be successfully handled through transfer learning

306 techniques but would require some labelled data from a new domain to facilitate such a feedback

307 training loop. The choice of data imputation technique could be another source of minor but

308 systematic error. The discriminant power of the model to detect HFpEF is lower in test subsets

309 where the missing data rate is higher and HFpEF cases are a small proportion of the overall

310 number. Finally, the applicability of our model in patients with HFpEF who have never required

311 hospital evaluation or admission is unknown. However, a strength of our approach is that a

312 dedicated specialist assessment for HF is not required to assess the probability of HFpEF among

313 patients undergoing general hospital evaluation (e.g. non-cardiological), even in the absence of

314 commonly used diagnostic data such as NTproBNP levels. The lack of independent validation is  
315 a limitation of this study. Evaluation of the derived model's performance in independent datasets  
316 from other centres and in community-based datasets will be informative in future studies.

317 Although we compared performance of the model with the H<sub>2</sub>FPEF score,[16] due to its stated  
318 aim of estimating the likelihood that HFpEF among patients with unexplained dyspnoea to guide  
319 further testing, we did not compare performance to the HFA-PEFF algorithm[37] which is a  
320 multi-step diagnostic algorithm. Furthermore, the comparison of our algorithm's performance  
321 with the H<sub>2</sub>FPEF should be confirmed in a separate validation cohort.

322

### 323 *Conclusion*

324 In this study, we have developed a rapid and automated data-driven approach that is effective at  
325 identifying patients from EHR who are likely to have HFpEF. This algorithm affords significant  
326 potential to rapidly identify patients for more detailed analyses and/or potential inclusion in  
327 clinical trials. The approach that we report could in principle be readily applied to other diseases  
328 and conditions that are similarly difficult to diagnose.

329

330 **Supplemental Materials.** The supplementary digital content is provided to support the findings  
331 of this study.

332

### 333 **Contributors**

334 Study design: NF, KO, RD, AMS;

335 Data collection: NF, JK, JO, AM;  
336 Data modeling: NF, DB, RD;  
337 Data analysis: NF, KO;  
338 Clinical validation: KO, RZ, DB2, AN;  
339 Result interpretation and writing the paper: All authors.  
340 Funding: AMS.  
341 Supervision: RD and AMS.

### 342 **Acknowledgements**

343 We thank Ahmed Mahmud and Joe Omigie for their invaluable advice during the data  
344 collection phase and Norman Catibog and Thiago Fonseca for sharing their knowledge in  
345 echocardiography. This work was supported by the British Heart Foundation (RE/18/2/34213;  
346 CH/1999001/11735); the NIHR Biomedical Research Centres at Guy's & St Thomas' NHS  
347 Foundation Trust (IS-BRC-1215-20006) and South London and Maudsley NHS Foundation  
348 Trust (IS-BRC-1215-20018), both with King's College London. KOG is supported by a Medical  
349 Research Council Clinical Training Fellowship (MR/R017751/1). DMB is funded by a UKRI  
350 Innovation Fellowship as part of Health Data Research UK MR/S00310X/1  
351 (<https://www.hdr.uk.ac.uk>). The views expressed are those of the authors and not necessarily  
352 those of NIHR or the Department of Health and Social Care. The funders had no role in study  
353 design, data collection and analysis, decision to publish, or preparation of the manuscript.

354

### 355 **Conflicts of Interest**

356 The authors have no conflicts of interest to declare.

357

358 **Data Sharing** The raw data used in this research are not openly available.

359

360

361 **References**

- 362
- 363 1. Owan TE, Hodge DO, Herges RM, Jacobsen SJ, Roger VL, Redfield MM. Trends in prevalence and  
364 outcome of heart failure with preserved ejection fraction. *N Engl J Med*. 2006;**355**:251-259.
  - 365 2. Bursi F, Weston SA, Redfield MM, Jacobsen SJ, Pakhomov S, Nkomo VT, Meverden RA, Roger VL.  
366 Systolic and diastolic heart failure in the community. *JAMA*. 2006;**296**:2209-2216.
  - 367 3. Zile MR, Baicu CF, Gaasch WH. Diastolic heart failure--abnormalities in active relaxation and  
368 passive stiffness of the left ventricle. *N Engl J Med*. 2004;**350**:1953-1959.
  - 369 4. Zile MR, Baicu CF, Ikonomidis JS, Stroud RE, Nietert PJ, Bradshaw AD, Slater R, Palmer BM, Van  
370 Buren P, Meyer M, Redfield MM, Bull DA, Granzier HL, LeWinter MM. Myocardial stiffness in patients  
371 with heart failure and a preserved ejection fraction: contributions of collagen and titin. *Circulation*.  
372 2015;**131**:1247-1259.
  - 373 5. Kawaguchi M, Hay I, Fetics B, Kass DA. Combined ventricular systolic and arterial stiffening in  
374 patients with heart failure and preserved ejection fraction: implications for systolic and diastolic reserve  
375 limitations. *Circulation*. 2003;**107**:714-720.
  - 376 6. Sunagawa K, Maughan WL, Burkhoff D, Sagawa K. Left ventricular interaction with arterial load  
377 studied in isolated canine ventricle. *Am J Physiol*. 1983;**245**:H773-780.
  - 378 7. Leite-Moreira AF, Correia-Pinto J. Load as an acute determinant of end-diastolic pressure-  
379 volume relation. *Am J Physiol Heart Circ Physiol*. 2001;**280**:H51-59.
  - 380 8. Leite-Moreira AF, Correia-Pinto J, Gillebert TC. Afterload induced changes in myocardial  
381 relaxation: a mechanism for diastolic dysfunction. *Cardiovasc Res*. 1999;**43**:344-353.
  - 382 9. Borlaug BA, Melenovsky V, Redfield MM, Kessler K, Chang HJ, Abraham TP, Kass DA. Impact of  
383 arterial load and loading sequence on left ventricular tissue velocities in humans. *J Am Coll Cardiol*.  
384 2007;**50**:1570-1577.
  - 385 10. Reddy YNV, Andersen MJ, Obokata M, Koeppe KE, Kane GC, Melenovsky V, Olson TP, Borlaug BA.  
386 Arterial Stiffening With Exercise in Patients With Heart Failure and Preserved Ejection Fraction. *J Am Coll*  
387 *Cardiol*. 2017;**70**:136-148.
  - 388 11. Yusuf S, Pfeffer MA, Swedberg K, Granger CB, Held P, McMurray JJ, Michelson EL, Olofsson B,  
389 Ostergren J, Investigators C, Committees. Effects of candesartan in patients with chronic heart failure  
390 and preserved left-ventricular ejection fraction: the CHARM-Preserved Trial. *Lancet*. 2003;**362**:777-781.
  - 391 12. Solomon SD, McMurray JJV, Anand IS, Ge J, Lam CSP, Maggioni AP, Martinez F, Packer M, Pfeffer  
392 MA, Pieske B, Redfield MM, Rouleau JL, van Veldhuisen DJ, Zannad F, Zile MR, Desai AS, Claggett B,  
393 Jhund PS, Boytsov SA, Comin-Colet J, Cleland J, Dungen HD, Goncalvesova E, Katova T, Kerr Saraiva JF,  
394 Lelonek M, Merkely B, Senni M, Shah SJ, Zhou J, Rizkala AR, Gong J, Shi VC, Lefkowitz MP, Investigators  
395 P-H, Committees. Angiotensin-Nepriylisin Inhibition in Heart Failure with Preserved Ejection Fraction. *N*  
396 *Engl J Med*. 2019;**381**:1609-1620.
  - 397 13. Pitt B, Pfeffer MA, Assmann SF, Boineau R, Anand IS, Claggett B, Clausell N, Desai AS, Diaz R, Fleg  
398 JL, Gordeev I, Harty B, Heitner JF, Kenwood CT, Lewis EF, O'Meara E, Probstfield JL, Shaburishvili T, Shah  
399 SJ, Solomon SD, Sweitzer NK, Yang S, McKinlay SM, Investigators T. Spironolactone for heart failure with  
400 preserved ejection fraction. *N Engl J Med*. 2014;**370**:1383-1392.
  - 401 14. Shah SJ, Katz DH, Selvaraj S, Burke MA, Yancy CW, Gheorghide M, Bonow RO, Huang CC, Deo  
402 RC. Phenomapping for novel classification of heart failure with preserved ejection fraction. *Circulation*.  
403 2015;**131**:269-279.

- 404 15. Shah SJ, Kitzman DW, Borlaug BA, van Heerebeek L, Zile MR, Kass DA, Paulus WJ. Phenotype-  
405 Specific Treatment of Heart Failure With Preserved Ejection Fraction: A Multiorgan Roadmap.  
406 *Circulation*. 2016;**134**:73-90.
- 407 16. Reddy YNV, Carter RE, Obokata M, Redfield MM, Borlaug BA. A Simple, Evidence-Based  
408 Approach to Help Guide Diagnosis of Heart Failure With Preserved Ejection Fraction. *Circulation*.  
409 2018;**138**:861-870.
- 410 17. Huusko J, Purmonen T, Toppila I, Lassenius M, Ukkonen H. Real-world clinical diagnostics of  
411 heart failure patients with reduced or preserved ejection fraction. *ESC Heart Fail*. 2020;**7**:1039-1048.
- 412 18. Ponikowski P, Voors AA, Anker SD, Bueno H, Cleland JG, Coats AJ, Falk V, Gonzalez-Juanatey JR,  
413 Harjola VP, Jankowska EA, Jessup M, Linde C, Nihoyannopoulos P, Parissis JT, Pieske B, Riley JP, Rosano  
414 GM, Ruilope LM, Ruschitzka F, Rutten FH, van der Meer P, Authors/Task Force M. 2016 ESC Guidelines  
415 for the diagnosis and treatment of acute and chronic heart failure: The Task Force for the diagnosis and  
416 treatment of acute and chronic heart failure of the European Society of Cardiology (ESC) Developed with  
417 the special contribution of the Heart Failure Association (HFA) of the ESC. *Eur Heart J*. 2016;**37**:2129-  
418 2200.
- 419 19. Wu H, Toti G, Morley KI, Ibrahim ZM, Folarin A, Jackson R, Kartoglu I, Agrawal A, Stringer C, Gale  
420 D, Gorrell G, Roberts A, Broadbent M, Stewart R, Dobson RJB. SemEHR: A general-purpose semantic  
421 search system to surface semantic data from clinical notes for tailored care, trial recruitment, and  
422 clinical research. *J Am Med Inform Assoc*. 2018;**25**:530-537.
- 423 20. Jackson R, Kartoglu I, Stringer C, Gorrell G, Roberts A, Song X, Wu H, Agrawal A, Lui K, Groza T,  
424 Lewsley D, Northwood D, Folarin A, Stewart R, Dobson R. CogStack - experiences of deploying integrated  
425 information retrieval and extraction services in a large National Health Service Foundation Trust  
426 hospital. *BMC Med Inform Decis Mak*. 2018;**18**:47.
- 427 21. Kraljevic Z BD, Mascio A, Roguski L, Folarin A, Roberts A, Bendayan R, Dobson R. MedCAT -  
428 medical concept annotation tool. 2019.
- 429 22. Kraljevic Z ST, Shek A, Roguski L, Noor K, Bean D, Mascio A, Zhu L, Folarin AA, Roberts A,  
430 Bendayan R, Richardson MP, Stewart R, Shah AD, Wong WK, Ibrahim Z, Teo JT, Dobson RJB. Multi-  
431 domain Clinical Natural Language Processing with MedCAT: the Medical Concept Annotation Toolkit.  
432 2020.
- 433 23. (MD) B. UMLS Reference Manual 2009.
- 434 24. Bean DM, Teo J, Wu H, Oliveira R, Patel R, Bendayan R, Shah AM, Dobson RJB, Scott PA.  
435 Semantic computational analysis of anticoagulation use in atrial fibrillation from real world data. *PLoS*  
436 *One*. 2019;**14**:e0225625.
- 437 25. Wharton G, Steeds R, Allen J, Phillips H, Jones R, Kanagala P, Lloyd G, Masani N, Mathew T,  
438 Oxborough D, Rana B, Sandoval J, Wheeler R, O'Gallagher K, Sharma V. A minimum dataset for a  
439 standard adult transthoracic echocardiogram: a guideline protocol from the British Society of  
440 Echocardiography. *Echo Res Pract*. 2015;**2**:G9-G24.
- 441 26. Lang RM, Badano LP, Mor-Avi V, Afalalo J, Armstrong A, Ernande L, Flachskampf FA, Foster E,  
442 Goldstein SA, Kuznetsova T, Lancellotti P, Muraru D, Picard MH, Rietzschel ER, Rudski L, Spencer KT,  
443 Tsang W, Voigt JU. Recommendations for cardiac chamber quantification by echocardiography in adults:  
444 an update from the American Society of Echocardiography and the European Association of  
445 Cardiovascular Imaging. *Eur Heart J Cardiovasc Imaging*. 2015;**16**:233-270.
- 446 27. Bielinski SJ, Pathak J, Carrell DS, Takahashi PY, Olson JE, Larson NB, Liu H, Sohn S, Wells QS,  
447 Denny JC, Rasmussen-Torvik LJ, Pacheco JA, Jackson KL, Lesnick TG, Gullerud RE, Decker PA, Pereira NL,  
448 Ryu E, Dart RA, Peissig P, Linneman JG, Jarvik GP, Larson EB, Bock JA, Tromp GC, de Andrade M, Roger  
449 VL. A Robust e-Epidemiology Tool in Phenotyping Heart Failure with Differentiation for Preserved and  
450 Reduced Ejection Fraction: the Electronic Medical Records and Genomics (eMERGE) Network. *J*  
451 *Cardiovasc Transl Res*. 2015;**8**:475-483.

- 452 28. Major V, Surkis A, Aphinyanaphongs Y. Utility of General and Specific Word Embeddings for  
453 Classifying Translational Stages of Research. *AMIA Annu Symp Proc.* 2018;**2018**:1405-1414.
- 454 29. Chen T GC. XGBoost: A Scalable Tree Boosting System. *KDD '16: Proceedings of the 22nd ACM*  
455 *SIGKDD International Conference on Knowledge Discovery and Data Mining.* 2016:785-794.
- 456 30. Lundberg S S-IL. A unified approach to interpreting model predictions. *NIPS.* 2017.
- 457 31. Steyerberg EW, Harrell FE, Jr., Borsboom GJ, Eijkemans MJ, Vergouwe Y, Habbema JD. Internal  
458 validation of predictive models: efficiency of some procedures for logistic regression analysis. *J Clin*  
459 *Epidemiol.* 2001;**54**:774-781.
- 460 32. Donahue J HJ, Rodner E, Saenko K, Darrell T. Semi-supervised Domain Adaptation with Instance  
461 Constraints. *2013 IEEE Conference on Computer Vision and Pattern Recognition.* 3012:668-675.
- 462 33. Pfeffer MA, Shah AM, Borlaug BA. Heart Failure With Preserved Ejection Fraction In Perspective.  
463 *Circ Res.* 2019;**124**:1598-1617.
- 464 34. Parikh KS, Sharma K, Fiuzat M, Surks HK, George JT, Honarpour N, Depre C, Desvigne-Nickens P,  
465 Nkulikiyinka R, Lewis GD, Gombert-Maitland M, O'Connor CM, Stockbridge N, Califf RM, Konstam MA,  
466 Januzzi JL, Jr., Solomon SD, Borlaug BA, Shah SJ, Redfield MM, Felker GM. Heart Failure With Preserved  
467 Ejection Fraction Expert Panel Report: Current Controversies and Implications for Clinical Trials. *JACC*  
468 *Heart Fail.* 2018;**6**:619-632.
- 469 35. Ho JE, Zern EK, Wooster L, Bailey CS, Cunningham T, Eisman AS, Hardin KM, Zampierollo GA,  
470 Jarolim P, Pappagianopoulos PP, Malhotra R, Naylor M, Lewis GD. Differential Clinical Profiles, Exercise  
471 Responses, and Outcomes Associated With Existing HFpEF Definitions. *Circulation.* 2019;**140**:353-365.
- 472 36. Kaspar M, Fette G, Guder G, Seidlmayer L, Ertl M, Dietrich G, Greger H, Puppe F, Stork S.  
473 Underestimated prevalence of heart failure in hospital inpatients: a comparison of ICD codes and  
474 discharge letter information. *Clin Res Cardiol.* 2018;**107**:778-787.
- 475 37. Pieske B, Tschope C, de Boer RA, Fraser AG, Anker SD, Donal E, Edelmann F, Fu M, Guazzi M,  
476 Lam CSP, Lancellotti P, Melenovsky V, Morris DA, Nagel E, Pieske-Kraigher E, Ponikowski P, Solomon SD,  
477 Vasan RS, Rutten FH, Voors AA, Ruschitzka F, Paulus WJ, Seferovic P, Filippatos G. How to diagnose heart  
478 failure with preserved ejection fraction: the HFA-PEFF diagnostic algorithm: a consensus  
479 recommendation from the Heart Failure Association (HFA) of the European Society of Cardiology (ESC).  
480 *Eur J Heart Fail.* 2020;**22**:391-412.

481

482

483 **TABLES**

484 **Table 1. Baseline characteristics of patients.** The mean and SD (standard deviation) were  
485 obtained where the predictor distribution follows a normal distribution, whereas for predictors  
486 with a skewed distribution, the median and interquartile range (25<sup>th</sup>-75<sup>th</sup>) were used to report the  
487 statistics. To evaluate the distributional differences between cases and controls, the Mann-  
488 Whitney U test or the t test was acquired, where appropriate. Values in parentheses next to each  
489 predictor name indicate the data availability percentage.

490 \* Constraint-free assumption on our test sets resulted in predictors with either a singular value or  
491 a high proportion of missing values. In such cases, the computation of common statistics was not  
492 pragmatic and hence the NAN (Not A Number) value was reported, instead.

493 \*\* This predictor is only computed in the test cohort to enable the comparison with the H<sub>2</sub>FPEF  
494 score.

495 # 92.45% of HFpEF cases and controls had a BNP or pro-BNP level available.

496 Set I: patients with normal EF, no/normal BNP record, a HF ICD10 code and at least one HF and  
497 dyspnea reference in their EHR.

498 Set II: patients with normal EF, no/normal BNP record, no HF diagnostic code and at least one  
499 HF and dyspnea reference in their EHR.

500 Set III: patients with normal EF, no BNP record, no HF diagnostic code nor HF reference in the  
501 EHR, at least one report of their dyspnea in their EHR.

502 Set IV: patients with normal EF, raised BNP result with HF and dyspnea reference in their EHR  
503 but no HF diagnosis documented



504 (HF: heart failure, EF: ejection fraction, rEF: reduced EF, BNP: brain-natriuretic peptide test,  
505 EHR: electronic health record) .

506 The following ICD10 codes were used to define the comorbidities:

507 Hypertension: I10-I15, I60-I69; Diabetes mellitus: E10-E14; Atrial fibrillation: I48; Pulmonary  
508 hypertension: I27; Kidney Disease: N18, N28, I12-I15

509

510 **Table 2. Multivariable model performance using the 5-fold cross-validation in derivation**  
511 **dataset.**

512

513 **Table 3. Additive SHAP feature importance for each category of predictors in the**  
514 **combined signatures.**

515

516 **Table 4. Multivariable model performance in independent test cohort.** The 95% CI is  
517 reported using bootstrapping in a thousand of iterations.

518 \*: HFpEF annotation agreement between the two scoring systems using Cohen's kappa statistics  
519 (python 3, Sklearn v.0.22).

520 AUROC: area under receiver operative curve, AP: average precision, CI: confidence interval in  
521 bootstrapped samples

522

523 **FIGURES**

524 **Figure 1. Feature selection analysis.** Features were incrementally utilized for training the  
525 models to ensure a performance within  $\pm 2$  units of the AUROC and f1-macro metrics in 5-fold  
526 cross-validation setup. Blue: f1-macro, Red: AUROC

527 **Figure 2. Feature importance using SHAP analysis in combined signatures.** Denser  
528 distribution of red points at the positive quadrant of the plot is representative of higher values of  
529 a given predictor's contribution in characterizing the positive class distribution i.e. in  
530 characterizing HFpEF.

531 **Figure 3. Performance of base and aggregate models.** Panel A: Receiver Operating  
532 Characteristic curves for base models, aggregate model, and H<sub>2</sub>FPEF score. Panel B: Precision  
533 Recall curves for base models, aggregate model, and H<sub>2</sub>FPEF score. Panel C: Calibration curve  
534 for aggregate model. Panel D: Efficiency curve for aggregate model. Panel E: Aggregate model  
535 performance in the 4 test subsets

536

537

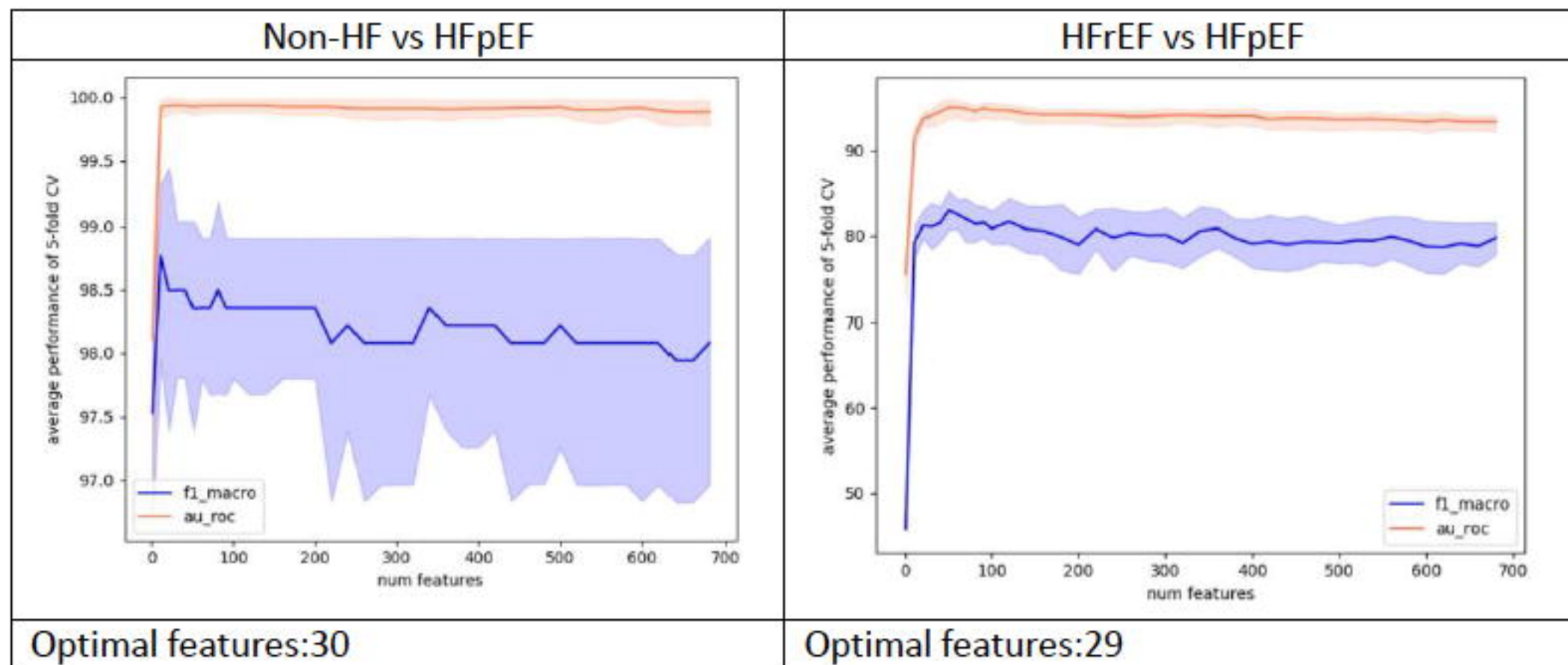


Figure 1. Feature selection analysis. Features were incrementally utilised for training the models to ensure a performance within  $\pm 2$  units of the AUROC and f1-macro metrics in 5-fold cross-validation setup. Blue: f1-macro, Red: AUROC

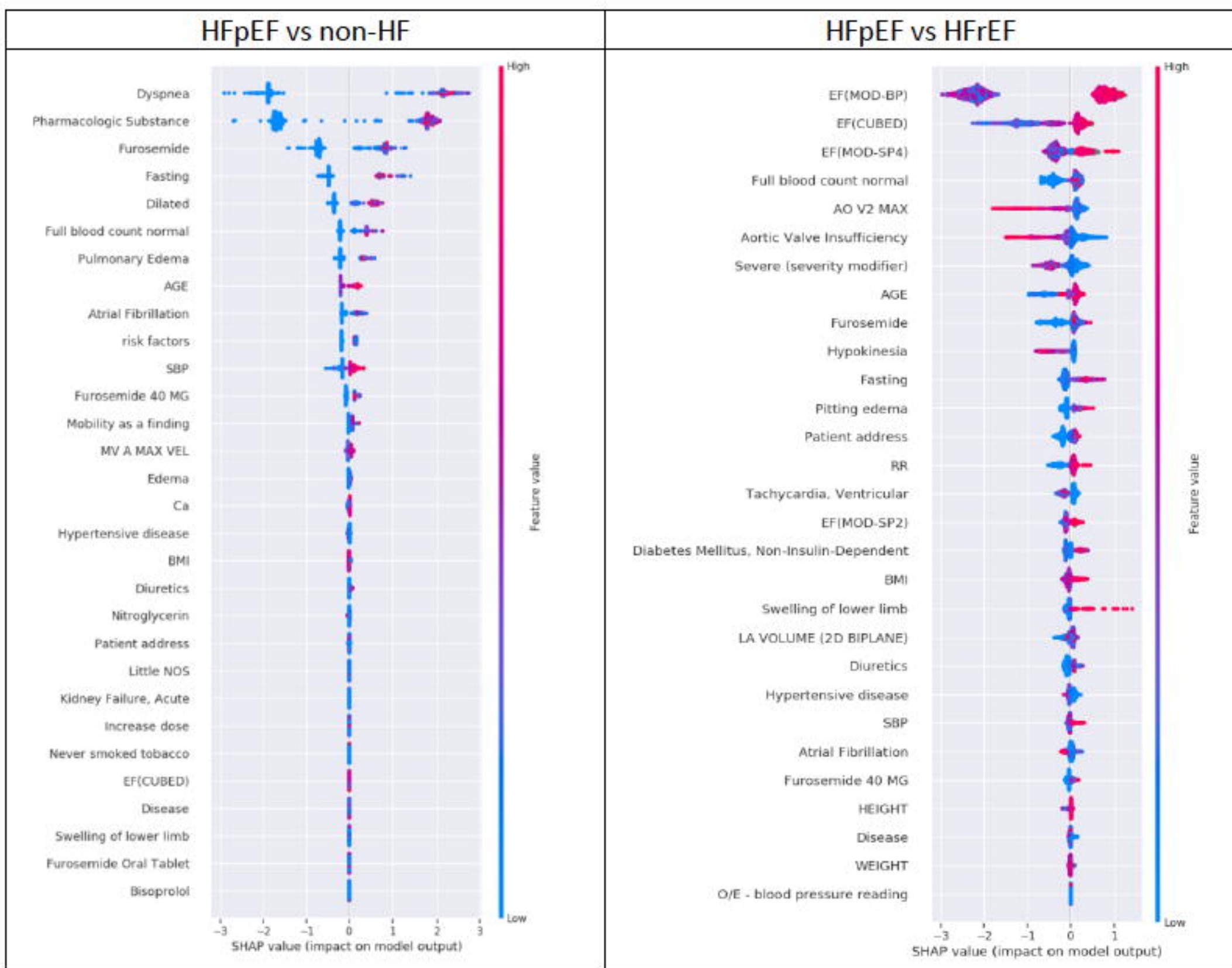
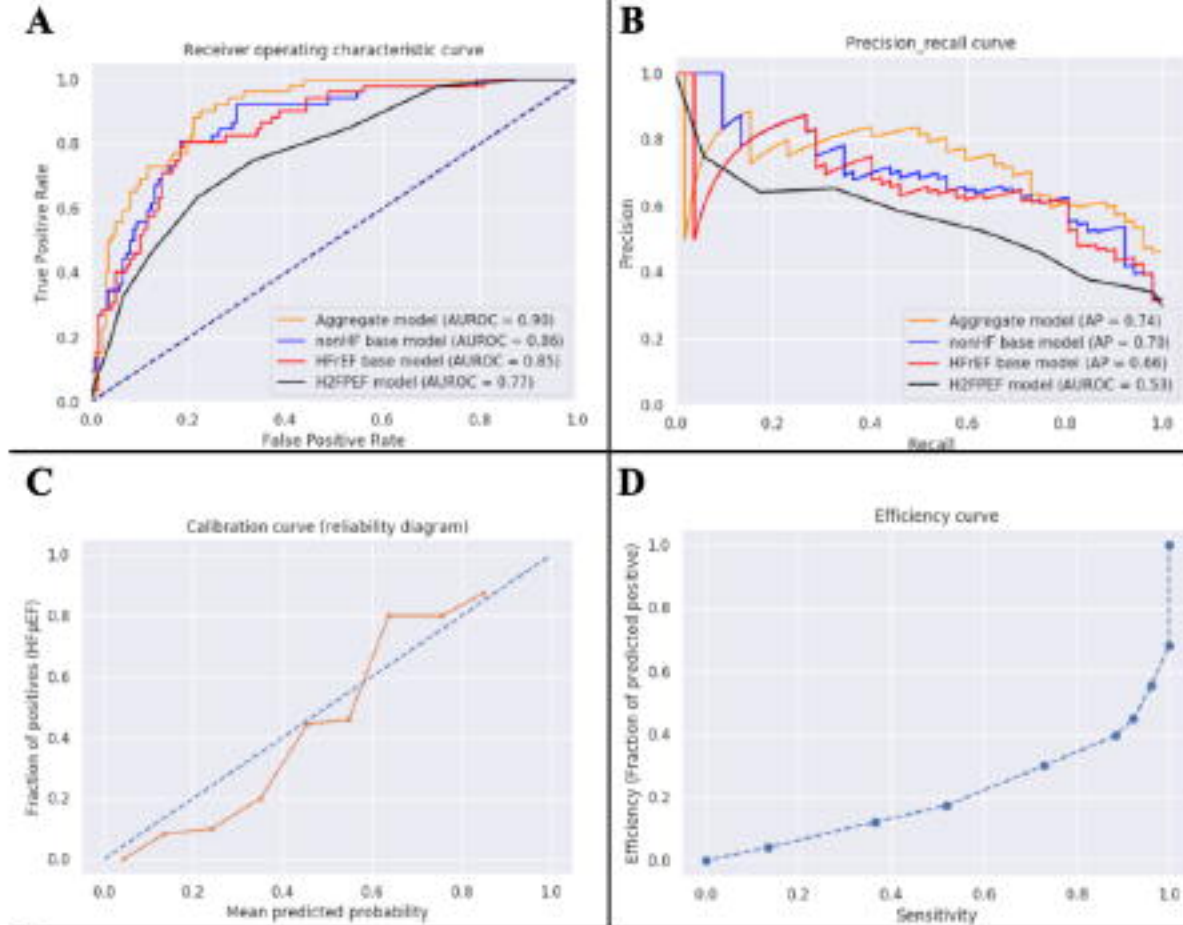


Figure 2. Feature importance using the SHAP analysis in combined signatures: denser distribution of red points at the positive quadrant of the plot is representative of higher values of a given predictor's contribution in characterizing the positive class distribution i.e. in characterizing HFpEF.



**E Aggregate model performance in the four test subsets**

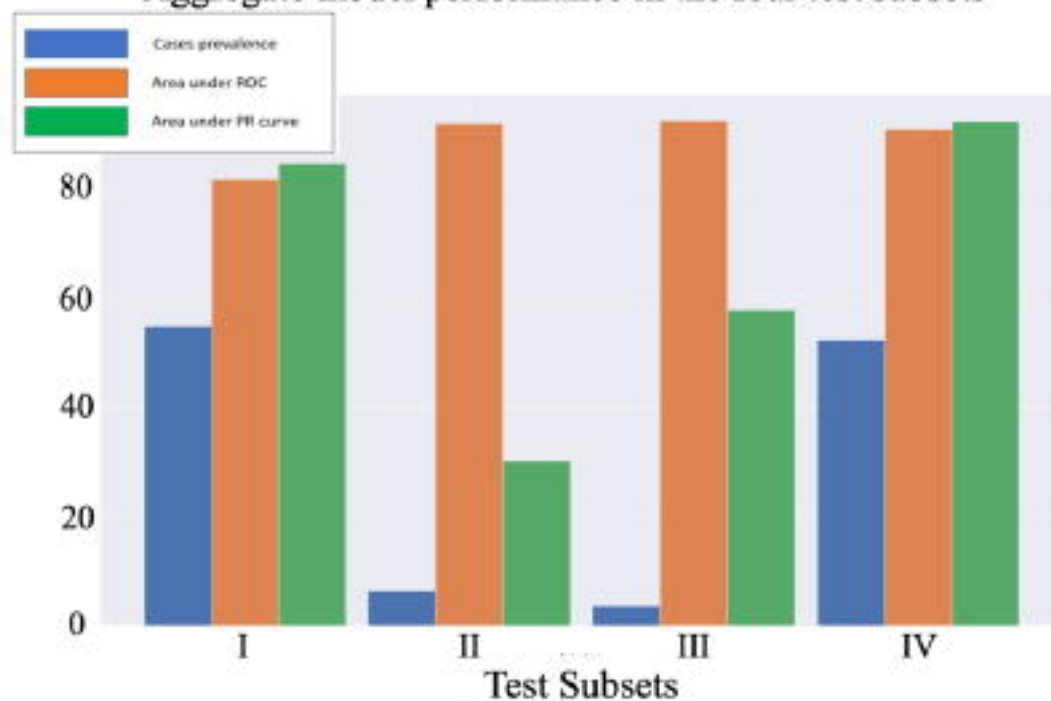


Figure 3. Performance of base and aggregate models. Panel A: Receiver Operating Characteristic curves for base models, aggregate model, and H2FPEF score. Panel B: Precision Recall curves for base models, aggregate model, and H2FPEF score. Panel C: Calibration curve for aggregate model. Panel D: Efficiency curve for aggregate model. Panel E: Aggregate model performance in the 4 test subsets



	Non-HF controls (n =194)	HF rEF controls (n=1258)	HFpEF cases (n =133)	P value cases vs controls	Test cohorts (n=269, HFpEF cases = 68)			
					Set I (n=61)	Set II (n=68)	Set III (n=71)	Set IV (n=69)
Female, % (100%)	48.5%	36.8%	54.9%	-	61.8%	67.1%	68.9%	61.2%
Age, y (100%)	54 ± 18	69 (22)	73 ± 12	<0.0001	66±13	56±15	55±15	61±13
Body mass index, kg/m <sup>2</sup> (76%)	28.35 ± 8.07	28.75 ± 7.34	34.06 ± 10.07	<0.0001	30.95 ± 8.15	32.18 ± 8.32	30.66 ± 7.58	31.67 ± 7.87
Hypertension, %	43.2%	81.6%	91.7%	-	83.8%	89.5%	67.6%	79.6%
Diabetes mellitus, %	20.1%	42%	54.1%	-	52.9%	31.6%	24.3%	34.7%
Atrial fibrillation, %	4.6%	47.6%	52.6%	-	50%	19.7%	6.7%	37.8%
Pulmonary hypertension, %	<1%	12.2%	25.6%	-	26.5%	7.9%	2.7%	11.2%
Kidney disease, %	6.7%	35.5%	46.6%	-	66.1%	21.1%	24.3%	25.5%
Antihypertensive drugs, n**	-	-	-	-	2(10)	0 (4)	0(4)	0(0)
NT-proBNP, pg/ml (#)	46 (53)	138 (1676)	4181 (3620)	-	873 (1359)	282 (181)	NAN*	781 (1258)
BNP, pg/ml (#)	54 (73)	76 (353)	1510 (4488)	-	NAN*	NAN*	NAN*	796 (656)
Creatinine, umol/l (99%)	82.8 ± 39.7	88.0 (34.0)	84.0 (28.0)	0.165	89.0 (40.0)	78.0 (25.5)	78.5 (24.0)	86.6 ± 19.6
Hemoglobin, g/dl (96%)	12.6 ± 2.1	13.3 (2.6)	13.1 ± 1.8	0.836	12.7 ± 2.0	12.8 ± 1.7	12.6 ± 2.0	12.9 ± 2.1
White cell count, 10 <sup>9</sup> /l (100%)	7.1 (4.33)	7.54 (3.99)	7.43 (3.76)	0.141	6.94 (3.57)	6.64 (3.4)	7.28 (4.43)	6.74 (3.16)

C-reactive protein, mg/l (96%)	6.5 ± 3.21	6.87 ± 3.12	7.4 (5.0)	0.254	6.93 ± 3.17	6.34 ± 3.01	6.62 ± 3.03	6.12 ± 3.11
Urea, mmol/l (99%)	5.73 ± 3.73	7.12 ± 4.43	6.4 (3.7)	0.687	6.85 (3.98)	5.3 (1.85)	4.65 (2.47)	5.95 (2.43)
Albumin, g/l (99%)	40.17 ± 6.98	41.13 ± 6.52	42.0 (3.0)	0.711	41.0 (6.0)	42.5 (4.25)	43.0 (6.0)	43.0 (3.0)
Sodium, mmol/l (99%)	138.34 ± 3.88	139.0 (4.0)	139.0 (3.0)	0.183	139.0 (3.25)	139.63 ± 2.52	139.34 ± 2.91	140.0 (3.0)
Potassium, mmol/l (99%)	4.57 ± 0.26	4.3 (0.6)	4.35 ± 0.58	0.720	4.2 (0.73)	4.28 ± 0.53	4.36 ± 0.52	4.31 ± 0.54
Calcium, mmol/l (98%)	2.28 (0.15)	2.29 (0.13)	2.31 ± 0.12	0.147	2.29 (0.17)	2.33 (0.14)	2.35 ± 0.13	2.34 ± 0.13
Systolic blood pressure, mmHg (63%)	132.93 ± 17.08	129.14 ± 21.9	139.54 ± 21.46	<0.0001	140.89 ± 25.87	136.78 ± 17.84	138.63 ± 23.21	138.12 ± 18.65
Diastolic blood pressure, mmHg (67%)	79.15 ± 11.76	73.79 ± 13.36	74.96 ± 13.45	<0.0001	73.16 ± 11.98	78.16 ± 13.09	80.68 ± 15.18	74.08 ± 13.73
Heart rate, beat/min (65%)	81.38 ± 14.28	73.69 ± 14.74	76.47 ± 16.72	0.0008	67.35 ± 9.63	75.31 ± 14.1	75.0 ± 8.8	74.64 ± 18.15
Oxygen saturation, % (52%)	98.1 ± 1.75	96.49 ± 5.05	96.22 ± 3.0	<0.0001	96.36 ± 2.87	96.69 ± 2.87	97.89 ± 1.62	96.13 ± 4.66
LV end diastolic volume, ml (23%)	155.0 ± 61.89	152.46 ± 53.52	106.7 ± 26.52	<0.0001	149.0 ± 16.97	124 ± NAN*	155.0 ± NAN *	110.83 ± 27.98
LV mass systolic, g (1%)	176.7 ± 53.0	265.2 ± 155.2	225.2 ± 74.5	<0.0001	118.9 ± nan*	NAN *	210.3 ± 158.0	111.7 ± 89.9
LV ejection fraction, % (100%)	60.4 ± 3.9	44.2 ± 11.5	58.0 ± 4.9	<0.0001	55.5 ± 2.1	60.5 ± 0.7	61.5 ± 2.1	55.3 ± 4.1
LV internal diameter at end diastole, cm/m <sup>2</sup> (59%)	2.46 ± 0.24	2.71 ± 0.5	2.46 ± 0.36	0.0002	2.36 ± 0.26	2.46 ± 0.28	2.32 ± 0.24	2.45 ± 0.35
LV stroke volume, ml (4%)	92.5 ± 34.78	65.52 ± 19.78	55.0 ± 4.36	<0.0001	82.0 ± 5.66	75.0 ± NAN *	93.0 ± NAN *	64.2 ± 19.7
LV outflow tract velocity time integral diameter, cm (20%)	2.13 ± 0.28	2.16 ± 0.24	2.17 ± 0.34	0.1126	2.03 ± 0.2	2.14 ± 0.24	2.13 ± 0.12	2.07 ± 0.23

LV end systolic volume, ml (22%)	62.5 ± 27.2	87.32 ± 42.81	42.67 ± 3.21	0.1708	66.5 ± 12.02	49.0 ± NAN *	62.0 ± NAN *	49.4 ± 10.45
LA systolic volume, ml (31%)	60.0 ± 19.52	86.47 ± 38.77	143.67 ± 70.44	< 0.0001	120.5 ± 28.99	112.0 ± NAN *	69.0 ± NAN *	70.33 ± 21.4
TR max PG, mmHg (80%)	26.7 ± 10.2	29.4 ± 11.2	34.16 ± 12.9	< 0.0001	37.2 ± 13.8	24.9 ± 10.18	26.3 ± 8.3	32.3 ± 11.4
E/e' lateral ratio (50%)	7.26 ± 3.11	10.70 ± 6.07	11.59 ± 5.96	< 0.0001	13.54 ± 4.59	11.70 ± 5.16	9.35 ± 3.48	12.87 ± 7.27
E/e' septal ratio (50%)	9.75 ± 4.99	14.51 ± 7.71	14.37 ± 5.6	< 0.0001	16.83 ± 5.44	14.77 ± 6.2	11.5 ± 4.52	16.04 ± 7.53
RV V1 max, cm/sec (6%)	82.28 ± 17.7	71.98 ± 22.65	80.41 ± 19.01	0.0001	95.37 ± 12.52	80.87 ± 17.42	71.34 ± 9.96	77.79 ± 18.88
RV V1 mean, cm/sec (4%)	47.04 ± 5.91	47.38 ± 14.19	52.03 ± 11.62	0.0048	53.5 ± 10.47	59.05 ± 5.18	48.65 ± 6.76	50.6 ± 9.53
Mitral valve E/A ratio, (84%)	1.08 ± 0.42	1.37 ± 0.07	1.13 ± 0.61	< 0.0001	1.23 ± 0.76	0.9 ± 0.34	0.98 ± 0.36	1.04 ± 0.4
Mitral regurgitation max velocity, cm/sec (10%)	483.46 ± 68.77	495.48 ± 88.56	502.84 ± 93.06	0.021	592.08 ± 98.05	553.39 ± 31.3	NAN	505.68 ± 119.98
Tricuspid regurgitation max velocity, cm/sec (80%)	233.23 ± 31.74	264.47 ± 56.62	274.1 ± 56.34	< 0.0001	310.38 ± 61.48	254.82 ± 53.76	239.73 ± 41.81	277.2 ± 50.81

**Table 1. Baseline characteristics of patients.** The mean and SD (standard deviation) were obtained where the predictor distribution follows a normal distribution, whereas for predictors with a skewed distribution, the median and interquartile range (25<sup>th</sup>-75<sup>th</sup>) were used to report the statistics. To evaluate the distributional differences between cases and controls, the Mann-Whitney U test or the t test was acquired, where appropriate. Values in parentheses next to each predictor name indicate the data availability percentage.

\* *Constraint-free assumption on our test sets resulted in predictors with either a singular value or a high proportion of missing values. In such cases, the computation of common statistics was not pragmatic and hence the NAN value (Not a Number) was reported, instead.*

\*\* *This predictor is only computed in the test cohort to enable the comparison with the H2FPEF score.*

# *92.45% of HFpEF cases and controls had a BNP or pro-BNP level available.*

*Set I: patients with normal EF, no/normal BNP record, a HF ICD10 code and at least one HF and dyspnea reference in their EHR.*



*Set II: patients with normal EF, no/normal BNP record, no HF diagnostic code and at least one HF and dyspnea reference in their EHR.*

*Set III: patients with normal EF, no BNP record, no HF diagnostic code nor HF reference in the EHR, at least one report of dyspnea in their EHR.*

*Set IV: patients with normal EF, raised BNP result with HF and dyspnea reference in their EHR but no HF diagnosis documented*

*(HF: heart failure, EF: ejection fraction, rEF: reduced EF, BNP: brain-natriuretic peptide test, EHR: electronic health record) .*

*The following ICD10 codes were used to define the comorbidities:*

*Hypertension: I10-I15, I60-I69; Diabetes mellitus: E10-E14; Atrial fibrillation: I48; Pulmonary hypertension: I27; Kidney Disease: N18, N28, I12-I15*

	Control set	f1_macro $\pm$ 95% CI	f1_weighted $\pm$ 95% CI	AUROC $\pm$ 95% CI
Structured Signature	Non-HF	84.05 $\pm$ 2.7	84.18 $\pm$ 2.7	92.04 $\pm$ 1.4
	HFrEF	75.75 $\pm$ 2.1	87.22 $\pm$ 1.42	90.31 $\pm$ 3.5
Unstructured Signature	Non-HF	<b>98.81 <math>\pm</math> 1.3</b>	<b>98.82 <math>\pm</math> 1.3</b>	99.7 $\pm$ 0.5
	HFrEF	78.59 $\pm$ 4.9	88.99 $\pm$ 2.1	94.38 $\pm$ 1.4
combined	Non-HF	98.57 $\pm$ 1.4	98.59 $\pm$ 1.4	<b>99.8 <math>\pm</math> 0.3</b>
signature	HFrEF	<b>83.03 <math>\pm</math> 2.8</b>	<b>90.91 <math>\pm</math> 1.6</b>	<b>95.67 <math>\pm</math> 2.0</b>

**Table 1. Multivariable model performance using the 5-fold cross-validation in derivation dataset.**

		Unstructured data	Structured data			
	Model	Symptoms	Echocardiography parameters	Vitals	Age & Sex	Lab results
Summed importance of grouped features	HFpEF vs non-HF	0.953	0.036	0.011	0.033	<0.001
	HFpEF vs HFrfEF	0.551	0.334	0.115	0.058	<0.001

**Table 1. Additive SHAP feature importance for each category of predictors in the combined signatures.**

	Performance metric	H2FPEF <sup>16</sup> , %	Combined non-HF signature, %	Combined HFrEF signature, %	Aggregate model score, %	Scoring Agreement * (HFrEF and non_HF)
Test set	AUROC (95% CI)	0.77	0.86 (± 0.002)	0.85 (± 0.001)	<b>0.90 (± 0.002)</b>	0.3
	AP	0.53	0.70	0.66	<b>0.74</b>	

**Table 1. Multivariable model performance in independent test cohort.** The 95% CI is reported using bootstrapping in a thousand of iterations.

\*: HFpEF annotation agreement between the two scoring systems using Cohen’s kappa statistics (python 3, Sklearn v.0.22).

AUROC: area under receiver operative curve, AP: average precision, CI: confidence interval in bootstrapped samples