

# 1 **Interpretable AI for beat-to-beat cardiac function assessment**

2  
3 David Ouyang<sup>1,\*</sup>, Bryan He<sup>2</sup>, Amirata Ghorbani<sup>3</sup>, Curt P. Langlotz<sup>4</sup>, Paul A. Heidenreich<sup>1</sup>, Robert  
4 A. Harrington<sup>1</sup>, David H. Liang<sup>1,3</sup>, Euan A. Ashley<sup>1,5,^</sup>, and James Y. Zou<sup>2,3,5,\*</sup>,^

- 5  
6 1. Department of Medicine, Stanford University  
7 2. Department of Computer Science, Stanford University  
8 3. Department of Electrical Engineering, Stanford University  
9 4. Department of Radiology, Stanford University  
10 5. Department of Biomedical Data Science, Stanford University

11 \* To whom correspondence should be addressed.

12 ^ Co-senior author.

13  
14 Correspondence: [ouyangd@stanford.edu](mailto:ouyangd@stanford.edu) and [jamesz@stanford.edu](mailto:jamesz@stanford.edu)

## 18 **Key Points**

- 19 ● Video based deep learning evaluation of cardiac ultrasound accurately identifies  
20 cardiomyopathy and predict ejection fraction, the most common metric of cardiac function.  
21 ● Using human tracings obtained during standard clinical workflow, deep learning semantic  
22 segmentation accurately labels the left ventricle frame-by-frame, including in frames  
23 without prior human annotation.  
24 ● Computational cardiac function analysis allows for beat-by-beat assessment of ejection  
25 fraction, which more accurately assesses cardiac function in patients with atrial fibrillation,  
26 arrhythmias, and heart rate variability.

27  
28  
29  
30  
31  
32  
33

34

## 35 **Abstract**

36 Accurate assessment of cardiac function is crucial for diagnosing cardiovascular disease<sup>1</sup>,  
37 screening for cardiotoxicity<sup>2,3</sup>, and deciding clinical management in patients with critical illness<sup>4</sup>.  
38 However human assessment of cardiac function focuses on a limited sampling of cardiac cycles  
39 and has significant interobserver variability despite years of training<sup>2,5,6</sup>. To overcome this  
40 challenge, we present the first beat-to-beat deep learning algorithm that surpasses human expert  
41 performance in the critical tasks of segmenting the left ventricle, estimating ejection fraction, and  
42 assessing cardiomyopathy. Trained on echocardiogram videos, our model accurately segments the  
43 left ventricle with a Dice Similarity Coefficient of 0.92, predicts ejection fraction with mean  
44 absolute error of 4.1%, and reliably classifies heart failure with reduced ejection fraction (AUC of  
45 0.97). Prospective evaluation with repeated human measurements confirms that our model has less  
46 variance than experts. By leveraging information across multiple cardiac cycles, our model can  
47 identify subtle changes in ejection fraction, is more reproducible than human evaluation, and lays  
48 the foundation for precise diagnosis of cardiovascular disease. As a new resource to promote  
49 further innovation, we also make publicly available one of the largest medical video dataset of  
50 over 10,000 annotated echocardiograms.

51

52

## 53 **Introduction**

54 Cardiac function is essential for maintaining normal systemic tissue perfusion with cardiac  
55 dysfunction manifesting as dyspnea, fatigue, exercise intolerance, fluid retention and  
56 mortality<sup>1,3,4,6-9</sup>. Impairment of cardiac function is labeled as “cardiomyopathy” or “heart failure”  
57 and is a leading cause of hospitalization in the United States and a growing global health issue<sup>1,10,11</sup>.  
58 A variety of methodologies have been used to quantify cardiac function and diagnose dysfunction.  
59 In particular, left ventricular ejection fraction (EF), the ratio of left ventricular end systolic and  
60 end diastolic volume, is one of the most important metrics of cardiac function, as it identifies  
61 patients who are eligible for life prolonging therapies<sup>8,12</sup>. However, there can be significant  
62 interobserver variability as well as inter-modality discordance based on methodology and  
63 modality<sup>2,5,6,12-15</sup>.

64

65 Human assessment of ejection fraction has variance in part due to common irregularity in the heart  
66 rate and the laborious nature of calculation limiting every beat quantification<sup>5,6</sup>. While the  
67 American Society of Echocardiography and the European Association of Cardiovascular Imaging  
68 guidelines recommend tracing and averaging up to 5 consecutive cardiac cycles if variation is  
69 identified, EF is often evaluated from tracings of only one representative beat or visually  
70 approximated if a tracing is deemed inaccurate<sup>6</sup>. This results in high variance and limited  
71 precision<sup>6,16</sup> with interobserver variation ranging from 7.6% to 13.9%<sup>5,13-16</sup>. This variation is  
72 observed despite substantial training by those reading the EF. More precise evaluation of cardiac  
73 function is necessary, as even patients with borderline reduction in EF have been shown to have  
74 significantly increased morbidity and mortality<sup>17-19</sup>.

75  
76 With rapid image acquisition, relatively low cost, and without ionizing radiation,  
77 echocardiography is the most widely used modality for cardiovascular imaging<sup>20,21</sup>. Being the most  
78 common first-line cardiovascular imaging modality, there is great interest in using deep learning  
79 techniques to determine ejection fraction<sup>22-24</sup>. Limitations in human interpretation, including  
80 laborious manual segmentation and inability to perform beat-to-beat quantification may be  
81 overcome by sophisticated automated approaches<sup>6,25,26</sup>. Recent advances in deep learning suggest  
82 that it can accurately and reproducibly identify human-identifiable phenotypes as well as  
83 characteristics unrecognized by human experts<sup>25,27-29</sup>.

84  
85 To overcome current limitations of human assessment of the left ventricular ejection fraction, we  
86 propose EchoNet-Dynamic, an end-to-end deep learning approach for left ventricular labeling and  
87 ejection fraction estimation from input echocardiogram videos alone. We first perform frame-level  
88 semantic segmentation of the left ventricle with weak supervision from prior clinical expert  
89 labeling. The segmentations are then combined with the native echocardiogram videos as input for  
90 a 3-dimensional (3D) convolutional neural network (CNN) with residual connections. This  
91 approach provides interpretable tracings of the ventricle, which facilitate human assessment and  
92 downstream analysis, while leveraging the 3D CNN to fully capture spatiotemporal patterns in the  
93 video<sup>6,30,31</sup>.

94

95

## 96 **Results**

97 EchoNet-Dynamic has three key components (Figure 1). First, we constructed a CNN model with  
98 atrous convolutions for frame-level semantic segmentation of the left ventricle. Atrous  
99 convolutions has been previously shown to perform well on non-medical imaging datasets<sup>30</sup>. The  
100 standard human clinical workflow for estimating ejection fraction requires manual segmentation  
101 of the left ventricle during end-systole and end-diastole. We generalize these labels in a weak  
102 supervision approach with atrous convolutions to generate frame-level semantic segmentation  
103 throughout the cardiac cycle in a 1:1 pairing with the original video. This automatic segmentation  
104 improves the robustness of our model and make it more interpretable to clinicians.

105  
106 Second, we trained a CNN model with residual connections and 3D spatiotemporal convolutions  
107 across frames to predict ejection fraction. Unlike prior 3D CNN architectures for medical imaging  
108 machine learning, our approach integrates spatial as well as temporal information with temporal  
109 variation across frames as the third dimension in our network convolutions<sup>25,31,32</sup>. Spatiotemporal  
110 convolutions, which incorporate spatial information in two dimensions as well as temporal  
111 information in the third dimension has been previously used in non-medical video classification  
112 tasks<sup>31,32</sup>, however has not been previously attempted on medical imaging given the relative  
113 scarcity of video medical imaging datasets nor used for regression tasks instead of classification  
114 tasks.

115  
116 Finally, we make video-level predictions of ejection fraction for beat-to-beat estimation of cardiac  
117 function. Each echocardiogram video typically includes multiple cardiac cycles, or beats, with  
118 each cycle being sufficient to produce a point estimate for ejection fraction. Given variance in  
119 cardiac function caused by changes in loading conditions as well as heart rate in a variety of cardiac  
120 conditions, it is recommended to perform ejection fraction estimation in up to 5 cardiac cycles,  
121 however this is not always done in clinical practice given the tedious and laborious nature of the  
122 calculation<sup>6,16</sup>. Our model identifies each cardiac cycle, generates a subsampled video-clip of 32  
123 frames, and averages clip-level estimates of EF as a form of test-time augmentation. Details of the  
124 model and hyperparameter search is further described in Methods, Supplementary Table 1, and  
125 Supplementary Figure 1.

126

127 EchoNet-Dynamic was developed using 10,030 apical-4-chamber echocardiograms obtained  
128 through the course of routine clinical practice at Stanford hospital. Each echocardiogram video  
129 corresponds to a unique patient during a unique visit and is representative of the variation in patient  
130 characteristics and image acquisition at the hospital. Table 1 contains the summary statistics of the  
131 patient population. These randomly selected patients have a range of ejection fractions  
132 representative of the patient population going through the echocardiography lab and the  
133 echocardiogram videos were split 7,465, 1,277, and 1,288 patients respectively for the training,  
134 validation, and test sets.

135

136 We worked with Stanford University and Hospitals to release our full dataset of 10,030 de-  
137 identified echocardiogram videos as a resource for the medical machine learning community for  
138 future comparison and validation of deep learning models. To the best of our knowledge, this is  
139 the largest labeled medical video dataset to be made publicly available and first large release of  
140 echocardiogram data with matched labels of human expert tracings, volume estimates, and left  
141 ventricular ejection fraction calculation. We expect this dataset to greatly facilitate new  
142 echocardiogram and medical video based machine learning work.

143

144 In a test dataset not previously seen during model training, model performance on individual  
145 subsampled video clips of approximately 1 second had a mean absolute error of 4.2% (95% CI  
146 4.0% - 4.3%), root mean squared error of 5.6% (5.7% - 5.8%) and  $R^2$  of 0.79 (95% CI 0.77 - 0.81)  
147 compared with the clinician report (Figure 2). Given that the model is agnostic to cardiac rhythm  
148 disturbances, including premature atrial contractions, premature ventricular contractions, and atrial  
149 fibrillation, we perform test time augmentation with beat-to-beat evaluation of ejection fraction.  
150 The final model with augmentation has improved performance with mean absolute error of 4.1%,  
151 root mean squared error of 5.3% and  $R^2$  of 0.81 (95% CI 0.78 - 0.82), which are within the range  
152 of typical measurement variation between different clinicians. We compared EchoNet-Dynamic's  
153 performance with that of several additional deep learning models that we trained on this dataset,  
154 and EchoNet-Dynamic is consistently more accurate, suggesting the power of its specific  
155 architecture (Supplementary Table 1).

156

157 EchoNet-Dynamic was compared against human measurements on 55 patients prospectively  
158 evaluated by two different sonographers on the same day. Each patient was independently  
159 evaluated for global longitudinal strain (GLS) and ejection fraction by multiple methods as well  
160 as our model for comparison (Figure 2D). EchoNet-Dynamic assessment of cardiac function had  
161 the least variance on repeat testing (median difference of 2.6%, SD=6.4) compared to EF obtained  
162 by Simpson's biplane method (median difference of 5.2%, SD=6.9,  $p < 0.001$  for non-inferiority),  
163 EF from Simpson's monoplane method (median difference of 4.6%, SD=7.3  $p < 0.001$  for non-  
164 inferiority), or GLS (median difference of 8.1%, SD=7.4%  $p < 0.001$  for non-inferiority). Of the  
165 initial 55 patients, 49 patients were also assessed with a different ultrasound system never seen  
166 during model training and EchoNet-Dynamic assessment had similar variance (median difference  
167 of 4.5%, SD=7.0,  $p < 0.001$  for non-inferiority for all comparisons with human measurements).

168  
169 EchoNet-Dynamic automatically generates segmentations of the left ventricle, which enables  
170 clinicians to better understand how it makes predictions. The segmentation is also useful because  
171 this provides a relevant point for human interjection in the workflow and physician oversight of  
172 the model in clinical practice. For the semantic segmentation task, the labels were 20,060 frame-  
173 level labels of the left ventricle obtained during the course of standard human clinical workflow  
174 during which expert human sonographers and echocardiographers manually label of the left  
175 ventricle during end-systole and end-diastole. Given the average video contains 2 labeled frames  
176 but 176 total frames, these weak labels were used to generate frame-level segmentations for the  
177 entire video (Figure 3). On the test dataset, the Dice Similarity Coefficient (DSC) for the end  
178 systolic tracing was 0.903 (95% CI 0.901 – 0.906) and the DSC for the end diastolic tracing was  
179 0.927 (95% CI 0.925 – 0.928). Despite being a frame-level, there was significant concordance in  
180 performance of end-systolic and end-diastolic semantic segmentation (Supplementary Figure 2).  
181 Example videos with semantic segmentation can be found in the Online Supplement.

182  
183 Variation in frame-to-frame model interpretation was seen in echocardiogram videos with  
184 arrhythmias and ectopy (Figure 4). In addition to correlation with irregularity in intervals between  
185 ventricular contractions, these videos were independently reviewed by clinical cardiologists and  
186 found to have atrial fibrillation, premature atrial contractions, and premature ventricular  
187 contractions. This highlights why it is important that EchoNet-Dynamic segments the ventricle

188 and estimates the EF for each every beat in the video and then aggregates across the beats. In  
189 particular, by aggregating across multiple beats, EchoNet-Dynamic significantly reduces variation  
190 compared to the common clinical practice of estimating EF from a single beat (Figure 4d).

191

## 192 **Discussion**

193 EchoNet-Dynamic is a new video deep learning technique that achieves state-of-the-art assessment  
194 of cardiac function. It uses expert human tracings for weak supervision of left ventricular  
195 segmentation, 3D spatiotemporal convolutions on video data, and beat-to-beat cumulative  
196 evaluation of EF across the entire video. EchoNet-Dynamic's performance in assessing EF is  
197 substantially better than prior deep learning attempts to assess EF<sup>22</sup>, and our model's variance is  
198 less than human expert measurements of cardiac function. EchoNet-Dynamic could potentially aid  
199 clinicians with more precise and reproducible assessment of cardiac function and detect subclinical  
200 change in ejection fraction beyond the precision of human readers. Furthermore, we release the  
201 largest annotated medical video dataset, which will stimulate future work on machine learning for  
202 cardiology.

203

204 EchoNet-Dynamic diverged the most from human estimation of EF in videos with arrhythmias  
205 and variation in heart rate. This variation is a feature of comparing EchoNet-Dynamic's beat-to-  
206 beat evaluation of EF across the video with our human evaluations of only one 'representative'  
207 beat. Choosing the representative beat can be subjective, contribute to human intra-observer  
208 variability, and less optimal compared to the guideline recommendation of averaging 5 consecutive  
209 beats. This workflow, is rarely done, in part due to the laborious and time intensive nature of the  
210 human tracing task. EchoNet-Dynamic greatly decreases the labor for cardiac function assessment  
211 with automating of the segmentation task and provide the opportunity for more frequent, rapid  
212 evaluation of cardiac function. Our end-to-end approach generates beat and clip level predictions  
213 of ejection fraction as well as segmentation of the left ventricle throughout the cardiac cycle for  
214 visual interpretation of the modeling results. In settings such as between chemotherapy sessions,  
215 after a heart transplant, and with the initiation of heart failure therapy, early detection of change in  
216 cardiac function significantly affect clinical care <sup>2,3</sup>.

217

218 Future studies will be required to ensure clinical applicability as well as generalizability in  
219 different clinical scenarios and health systems. With rapid expansion in the use of point of care  
220 ultrasound for evaluation of cardiac function by non-cardiologists, we aim to explore the feasibility  
221 and generalizability of our model with input videos are variable quality and acquisition expertise.  
222 Correlating the model performance with improved clinical outcomes and health system costs will  
223 also be required to determine potential impact. In addition to its application assessing left  
224 ventricular ejection fraction, the deep learning techniques applied in this study have considerable  
225 relevance to other types of medical video imaging data with temporal information, including  
226 cardiac magnetic resonance imaging, as well as other functional assessments using  
227 echocardiogram videos.

228  
229 These results represent an important step towards automated evaluation of cardiac function from  
230 echocardiogram videos through deep learning. EchoNet-Dynamic could augment current methods  
231 with improved precision, accuracy, and allow earlier detection of subclinical cardiac dysfunction,  
232 and the underlying dataset can be used to advance future work in deep learning for medical video  
233 imaging datasets and lay the groundwork for further applications of medical deep learning.

234

### 235 **Acknowledgements**

236 This work is supported by the Stanford Translational Research and Applied Medicine pilot grant,  
237 Stanford Cardiovascular Institute pilot grant, and a Stanford Artificial Intelligence in Imaging and  
238 Medicine Center seed grant. D.O. is supported by the American College of Cardiology Foundation  
239 / Merck Research Fellowship. B.H. is supported by the NSF Graduate Research Fellowship. A.G.  
240 is supported by the Stanford-Robert Bosch Graduate Fellowship in Science and Engineering. J.Z.  
241 is supported NSF CCF 1763191, NIH R21 MD012867-01, NIH P30AG059307 and by a Chan-  
242 Zuckerberg Biohub Fellowship.

243

### 244 **Author Contributions**

245 DO retrieved and quality controlled all videos and merged electronic medical record data. DO,  
246 BH, AG, JYZ developed and trained the deep learning algorithms, performed statistical tests, and  
247 created all the figures. DO, CPL, PAH, RAH coordinated public release of the deidentified



248 echocardiogram dataset. DO, PAH, DHL, EAA performed clinical evaluation of model  
249 performance. DO, BH, EAA, JYZ wrote the manuscript with all authors.

250

## 251 **Online Methods**

252

### 253 **Data Curation**

254 A standard full resting echocardiogram study consists of a series of 50-100 videos and still images  
255 visualizing the heart from different angles, locations, and image acquisition techniques (2D  
256 images, tissue Doppler images, color Doppler images, and others). In this dataset, one apical-4-  
257 chamber 2D gray-scale video is extracted from each study. Each video represents a unique  
258 individual as the dataset contains 10,025 echocardiography videos from 10,025 unique individuals  
259 who underwent echocardiography between 2016 and 2018 as part of clinical care at Stanford  
260 University Hospital. Images were acquired by skilled sonographers using iE33, Sonos, Acuson  
261 SC2000, Epiq 5G, or Epiq 7C ultrasound machines and processed images were stored in Philips  
262 Xcelera picture archiving and communication system. Video views were identified through  
263 implicit knowledge of view classification in the clinical database by identifying images and videos  
264 labeled with measurements done in the corresponding view.

265

266 The apical-4-chamber view video was identified by extracting the Digital Imaging and  
267 Communications In Medicine (DICOM) file linked to measurements of ventricular volume used  
268 to calculate the ejection fraction. Videos were spot checked for quality control, confirm view  
269 classification, and exclude videos with color Doppler. Each subsequent video was cropped and  
270 masked to remove text, ECG and respirometer information, and other information outside of the  
271 scanning sector. The resulting square images were either 600x600 or 768x768 pixels depending  
272 on the ultrasound machine and downsampled by cubic interpolation using OpenCV into  
273 standardized 112x112 pixel videos.

274

275 This research was approved by the Stanford University Institutional Review Board and data  
276 privacy review through a standardized workflow by the Center for Artificial Intelligence in  
277 Medicine and Imaging (AIMI) and the University Privacy Office. In addition to masking of text,  
278 ECG information, and extra data outside of the scanning sector in the video files as described

279 below, each DICOM file's pixel data was parsed out and saved as an AVI file to prevent any  
280 leakage of identifying information through public or private DICOM tags. Each video was  
281 subsequently manually reviewed by an employee of the Stanford Hospital with familiarity with  
282 imaging data to confirm the absence of any identifying information.

283

### 284 **Prospective Clinical Validation**

285 Prospective validation was performed by two senior sonographers with advanced cardiac  
286 certification and greater than 15 years experience each. For each patient, measurements of cardiac  
287 function was independently acquired, measured, and assessed by each sonographer on the same  
288 day. Every patient was scanned using Epiq 7C ultrasound machines, the standard instrument in the  
289 Stanford Echocardiography Lab, and a subset of patients were also rescanned by the same two  
290 sonographers using a GE Vivid 95E ultrasound machine. Tracing and measurement was done on  
291 a dedicated workstation after image acquisition. For comparison, the independently acquired  
292 apical-4-chamber videos were fed into the model and the variance in measurements assessed.

293

### 294 **EchoNet-Dynamic development and training**

295 Model building and training was done in Python on the PyTorch deep learning library. Semantic  
296 segmentation was performed using the Deeplabv3 architecture<sup>30</sup>. The segmentation model had a  
297 base architecture of 50-layer residual net and minimized a pixel level binary cross entropy loss.  
298 The model was initialized with random weights, and was trained using a stochastic gradient  
299 descent optimizer with a learning rate of 0.00001, momentum of 0.9, and batch size of 20 for 50  
300 epochs. Our model with spatiotemporal convolutions was initialized with pretrained weights from  
301 the Kinetics-400 dataset<sup>33</sup>. We tested three model architectures with variable integration of  
302 temporal convolutions and ultimately chose decomposed R2+1D spatiotemporal convolutions as  
303 the model with the best performance<sup>31,32</sup>. The models were trained to minimize the squared loss  
304 between the prediction and true ejection fraction using a stochastic gradient descent optimizer with  
305 an initial learning rate of 0.0001, momentum of 0.9, and batch size of 16 for 45 epochs. The  
306 learning rate was decayed by a factor of 0.1 every 15 epochs was used during model training.  
307 During training, clips of 32 frames were generated by sampled every other frame. To augment the  
308 clips, all frames were padded with 12 pixels, and a random crop of the 112x112 pixel size was

309 taken. For all models, the weights from the epoch with the lowest validation loss was selected for  
310 final testing.

### 311 312 **Test Time Augmentation with Beat-by-Beat Assessment**

313 There can be variation in the ejection fraction, end systolic volume, and end diastolic volumes  
314 during atrial fibrillation, and in the setting of premature atrial contractions, premature ventricular  
315 contractions, and other sources of ectopy. The clinical convention is to identify at least one  
316 representative cardiac cycle and use this representative cardiac cycle to perform measurements,  
317 although an average of the measurements of up to five cardiac cycles is recommended when there  
318 is significant ectopy or variation. For this reason, our final model used test time augmentation by  
319 providing individual estimates for each ventricular beat throughout the entire video and outputs  
320 the average prediction as the final model prediction. We use the segmentation model to identify  
321 the area of the left ventricle and threshold-based processing to identify ventricular contractions  
322 during each cardiac cycle. For beat, a subsampled clip centered around the ventricular contraction  
323 was obtained and used to produce a beat-by-beat estimate of EF. The mean ejection fraction of all  
324 ventricular contractions in the video was used as the final model prediction.

### 325 326 **Statistical Analysis**

327 Confidence intervals were computed using 10,000 bootstrapped samples and obtaining 95  
328 percentile ranges for each prediction. Chi-squared test and Student's t-test were used for statistical  
329 comparisons.

### 330 331 **Data Availability**

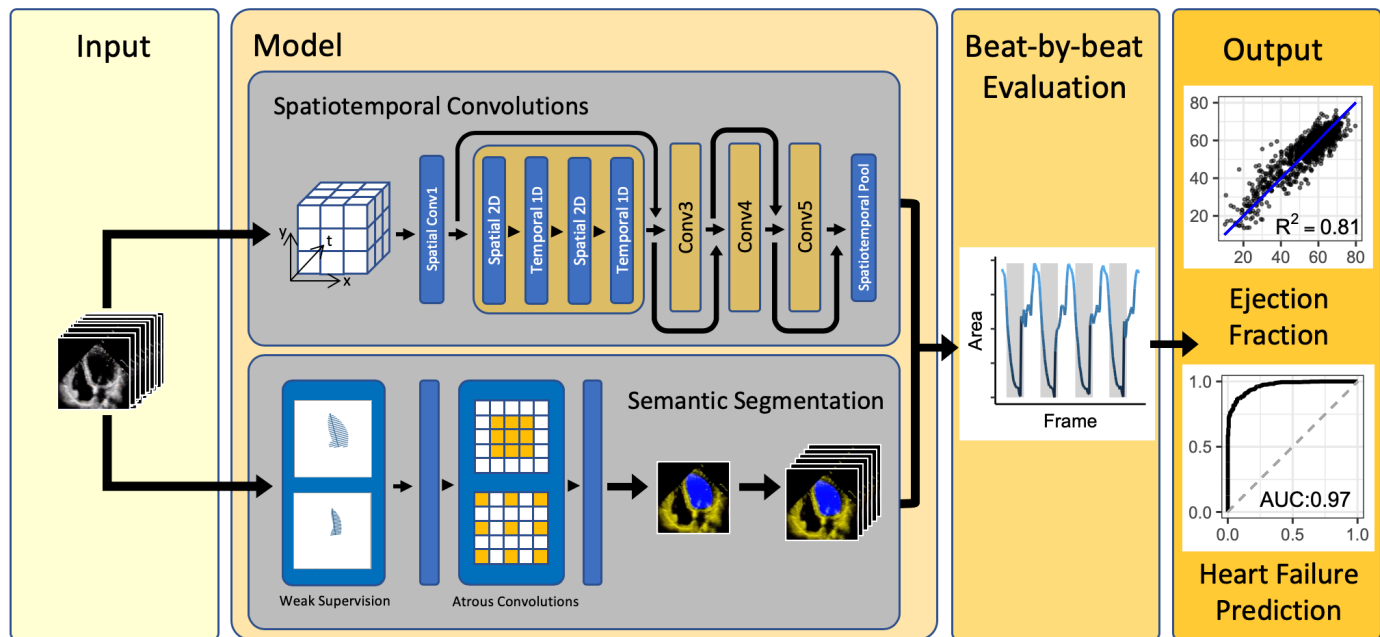
332 This data introduces the EchoNet-Dynamic Dataset, a publicly available dataset of deidentified  
333 echocardiogram videos, publicly available at: <https://douyang.github.io/EchoNetDynamic/>

### 334 335 **Code Availability**

336 The code is available at: <https://github.com/douyang/EchoNetDynamic/>

337  
338

339



340

341 Figure 1. EchoNet-Dynamic workflow. For each patient, EchoNet-Dynamic uses standard apical-  
342 4-chamber view echocardiogram video as input. The model first predicts ejection fraction for each  
343 cardiac cycle using 3D spatiotemporal convolutions with residual connections and generates  
344 frame-level semantic segmentations of the left ventricle using weak supervision from expert  
345 human tracings. These outputs are combined to create beat-by-beat predictions of ejection fraction  
346 and to predict the presence of heart failure with reduced ejection fraction.

347

348

349

350

351

352

353

354

355

356

357

358

359 Table 1. Summary statistics of patients in the dataset. Data obtained from visits to Stanford

360 Hospital between 2016 and 2018.

361

362

| Metric                         | Total       | Training    | Validation  | Test        |
|--------------------------------|-------------|-------------|-------------|-------------|
| Number of Patients             | 10,030      | 7,465       | 1,288       | 1,277       |
| Female, n (%)                  | 4,885 (49%) | 3,662 (49%) | 611 (47%)   | 612 (48%)   |
| Age, years (SD)                | 68 (21)     | 70 (22)     | 66 (18)     | 67 (17)     |
| Ejection Fraction, % (SD)      | 55.7 (12.5) | 55.7 (12.5) | 55.8 (12.3) | 55.3 (12.4) |
| End Systolic Volume, mL (SD)   | 43.3 (34.5) | 43.2 (36.1) | 43.3 (34.5) | 43.9 (36.0) |
| End Diastolic Volume, mL (SD)  | 91.0 (45.7) | 91.0 (46.0) | 91.0 (43.8) | 91.4 (46.0) |
| Heart Failure, n (%)           | 2,874 (29%) | 2,113 (28%) | 356 (28%)   | 405 (32%)   |
| Diabetes Mellitus, n (%)       | 2,018 (20%) | 1,474 (20%) | 275 (21%)   | 269 (21%)   |
| Hypercholesterolemia, n (%)    | 3,321 (33%) | 2,463 (33%) | 445 (35%)   | 413 (32%)   |
| Hypertension, n (%)            | 3,936 (39%) | 2,912 (39%) | 525 (41%)   | 499 (39%)   |
| Renal Disease, n (%)           | 2,004 (20%) | 1,475 (20%) | 249 (19%)   | 280 (22%)   |
| Coronary Artery Disease, n (%) | 2,290 (23%) | 1,674 (22%) | 302 (23%)   | 314 (25%)   |

363

364

365

366

367

368

369

370

371

372

373

374

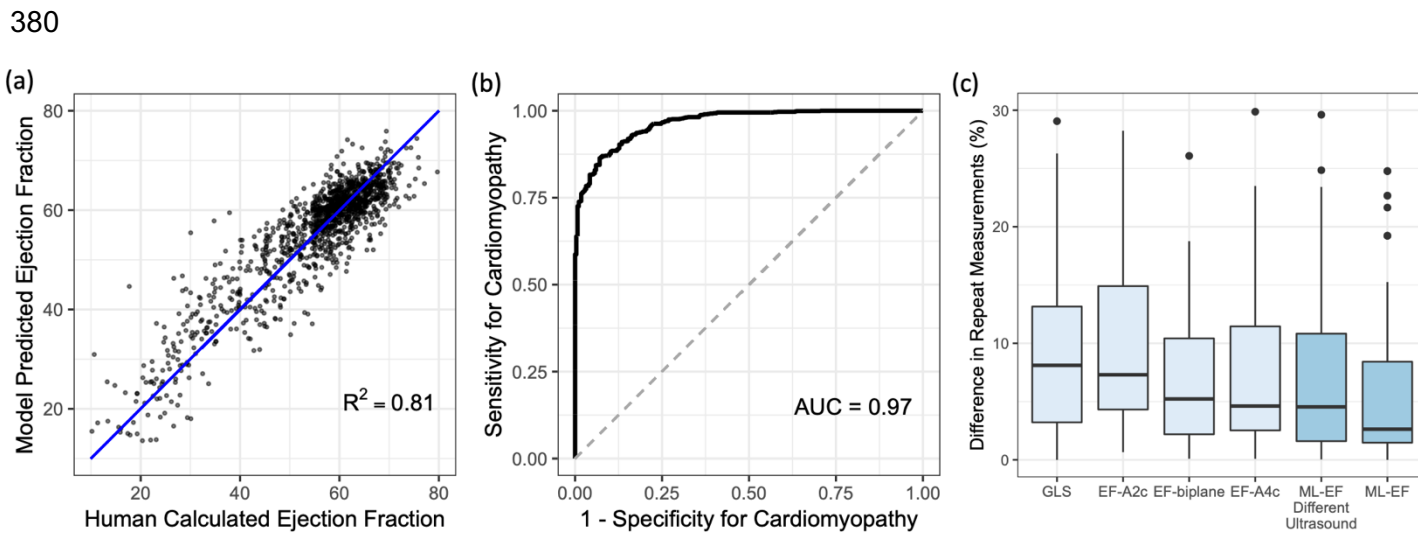
375

376

377

378

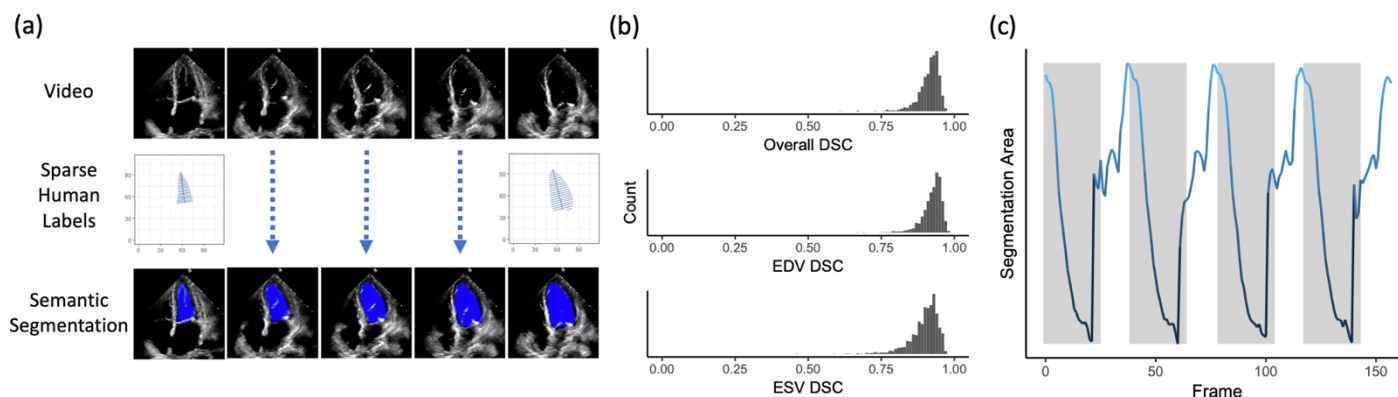
379



381  
382  
383 Figure 2. Model Performance. (a) EchoNet-Dynamic's predicted ejection fraction vs. reported  
384 ejection fraction. (b) Receiver operating characteristic curve for diagnosis of heart failure with  
385 reduced ejection fraction. (c) Variance of metrics of cardiac function on repeat measurement. The  
386 first four boxplots corresponds to variation by clinicians using different techniques, and the last  
387 two boxplots corresponds EchoNet-Dynamic's variance on input images from standard ultrasound  
388 machines and an ultrasound machine not previously seen by the model.

389  
390  
391  
392  
393  
394  
395  
396  
397  
398  
399  
400

401



402

403

404 Figure 3. Semantic Segmentation Performance. (a) Weak supervision with human expert tracings  
405 of the left ventricle at end-systole and end-diastole is used to train a semantic segmentation model  
406 with input video frames throughout the cardiac cycle. (b) Dice Similarity Coefficient (DSC) was  
407 calculated for each ESV/EDV frame. (c) The area of the left ventricle segmentation was used to  
408 identify heart rate and bin clips for beat-to-beat evaluation of ejection fraction.

409

410

411

412

413

414

415

416

417

418

419

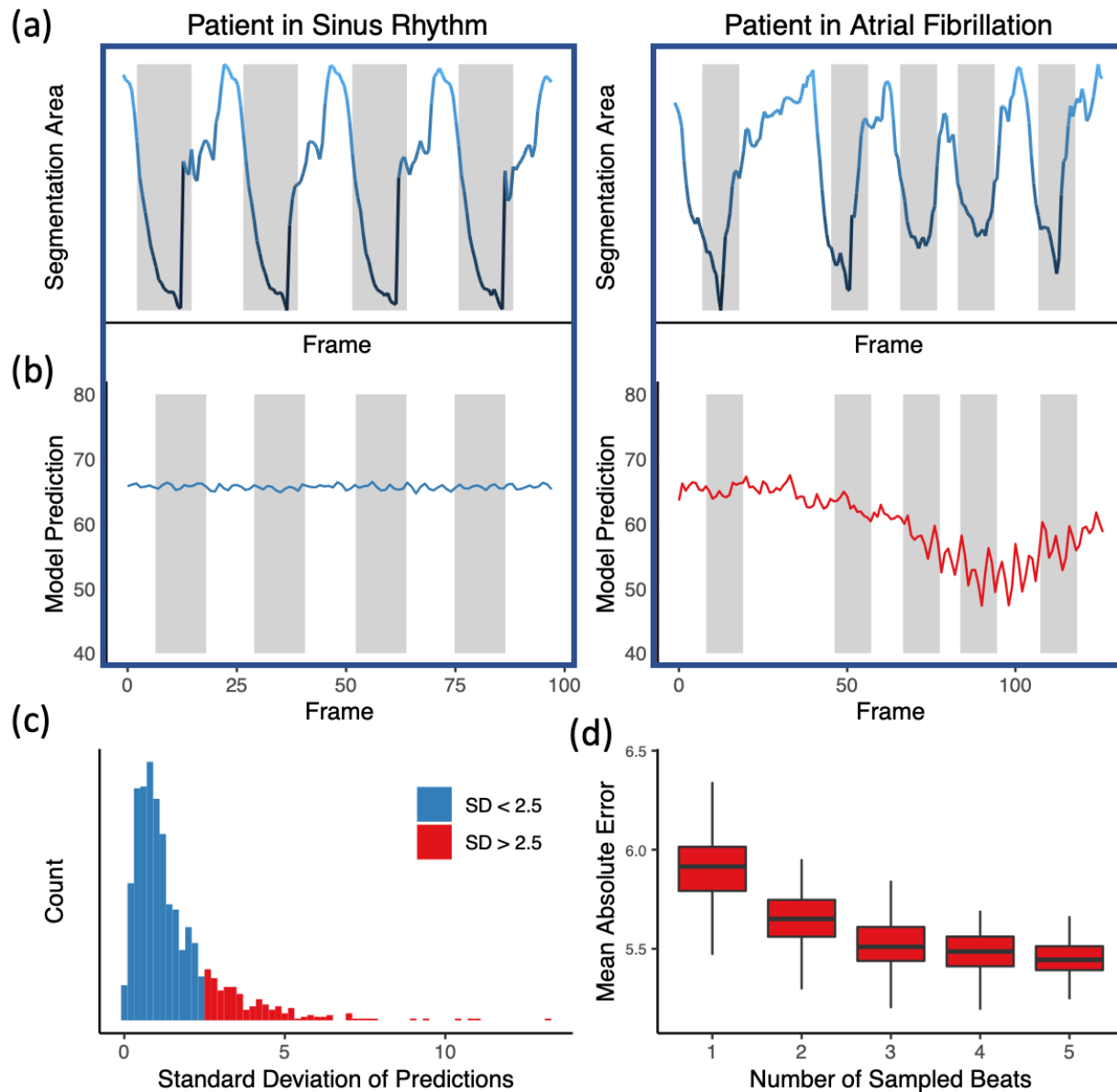
420

421

422

423

424



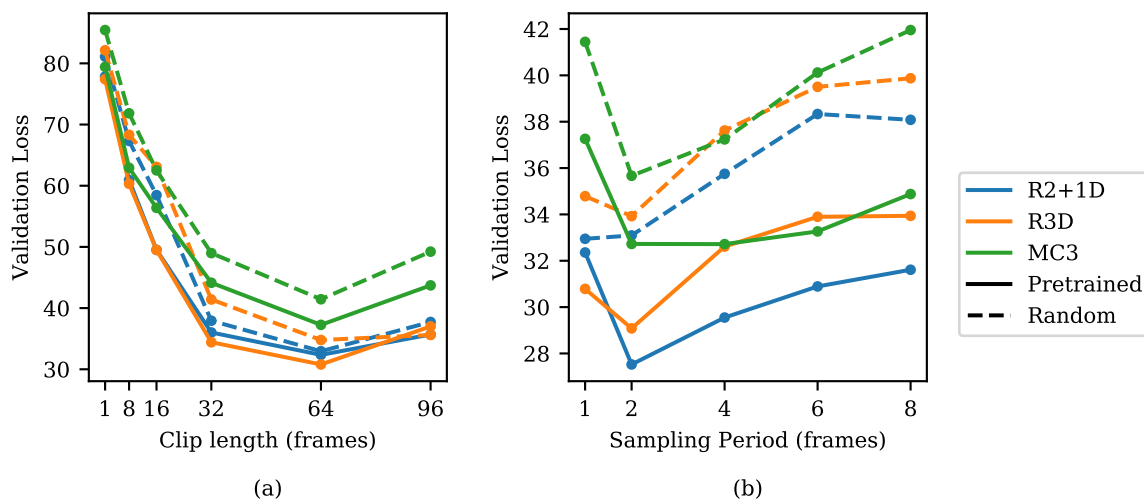
425  
426 Figure 4. Beat-to-beat evaluation of ejection fraction. (a) Atrial fibrillation and arrhythmias can be  
427 identified by significant variation in intervals between ventricular contractions. (b) Significant  
428 variation in left ventricle segmentation area was associated with higher variance in EF prediction.  
429 (c) Histogram of standard deviation of beat-to-beat evaluation of EF across all the test videos. (d)  
430 Assessing the effect of beat-to-beat based on the number of sampled beats averaged for prediction.  
431 Each boxplot represents 100 random samples of a certain number of beats and comparison with  
432 reported ejection fraction.

433

434



435



436

437 Supplementary Figure 1: Hyperparameter search for 3D Spatiotemporal Convolutions on video  
438 dataset to predict ejection fraction. Model architecture (R2+1D, R3D, and MC3), initialization  
439 (Kinetics-400 pretrained weights with solid line and random initial weights with dotted line), clip  
440 length (1, 8, 16, 32, 64, and 96), and sampling period (1, 2, 4, 8) were considered. (a) When varying  
441 clip lengths, performance is best at 64 frames (corresponding to 1.28 seconds), and starting from  
442 pretrained weights improves performance slightly across all models. (b) Varying sampling period  
443 with a length to approximately correspond to 64 frames prior to subsampling. Performance is best  
444 at a sampling period of 2.

445

446

447

448

449

450

451

452

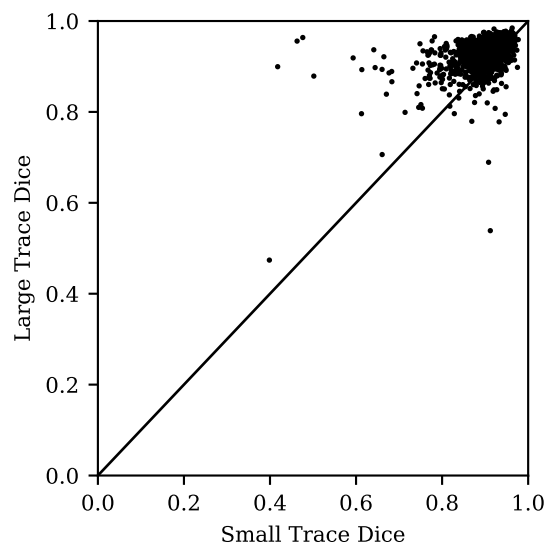
453

454

455

456

457



458

459 Supplementary Figure 2: Relationship between end systolic tracing Dice Similarity Coefficient  
460 and end diastolic tracing Dice Similarity Coefficient.

461

462

463 Supplementary Table 1: Model performance compared to three alternative deep learning models  
464 in assessing cardiac function.

465

| Model           | Evaluation      | Sampling Period | MAE  | RMSE | $R^2$ |
|-----------------|-----------------|-----------------|------|------|-------|
| EchoNet-Dynamic | Beat-by-beat    | 1 in 2          | 4.05 | 5.32 | 0.81  |
| R2+1D           | 32 frame sample | 1 in 2          | 4.22 | 5.56 | 0.79  |
| R3D             | 32 frame sample | 1 in 2          | 4.21 | 5.62 | 0.79  |
| MC3             | 32 frame sample | 1 in 2          | 4.54 | 5.97 | 0.77  |

466

467

468

469

470

471 **References**

- 472
- 473 1. Ziaieian, B. & Fonarow, G. C. Epidemiology and aetiology of heart failure. *Nat. Rev.*  
474 *Cardiol.* **13**, 368–378 (2016).
  - 475 2. Charbonnel, C. *et al.* Assessment of global longitudinal strain at low-dose anthracycline-  
476 based chemotherapy, for the prediction of subsequent cardiotoxicity. *Eur. Heart J.*  
477 *Cardiovasc. Imaging* **18**, 392–401 (2017).
  - 478 3. Shakir, D. K. & Rasul, K. I. Chemotherapy induced cardiomyopathy: pathogenesis,  
479 monitoring and management. *J. Clin. Med. Res.* **1**, 8–12 (2009).
  - 480 4. Dellinger, R. P. *et al.* Surviving Sepsis Campaign: International Guidelines for Management  
481 of Severe Sepsis and Septic Shock, 2012. *Intensive Care Med.* **39**, 165–228 (2013).
  - 482 5. Farsalinos, K. E. *et al.* Head-to-Head Comparison of Global Longitudinal Strain  
483 Measurements among Nine Different Vendors: The EACVI/ASE Inter-Vendor Comparison  
484 Study. *J. Am. Soc. Echocardiogr.* **28**, 1171–1181, e2 (2015).
  - 485 6. Lang, R. M. *et al.* Recommendations for cardiac chamber quantification by  
486 echocardiography in adults: an update from the American Society of Echocardiography and  
487 the European Association of Cardiovascular Imaging. *Eur. Heart J. Cardiovasc. Imaging*  
488 **16**, 233–270 (2015).
  - 489 7. Virchow, R. *Die Cellularpathologie in ihrer Begründung auf physiologische und*  
490 *pathologische Gewebelehre.* (Hirschwald, 1871).
  - 491 8. McMurray, J. J. *et al.* Task Force for the Diagnosis and Treatment of Acute and Chronic  
492 Heart Failure 2012 of the European Society of Cardiology; ESC Committee for Practice  
493 Guidelines. ESC guidelines for the diagnosis and treatment of acute and chronic heart  
494 failure 2012: The Task Force for the Diagnosis and Treatment of Acute and Chronic Heart

- 495 Failure 2012 of the European Society of Cardiology. Developed in collaboration with the  
496 Heart Failure Association (HFA) of the ESC. *Eur. J. Heart Fail.* **14**, 803–869 (2012).
- 497 9. Loehr, L. R., Rosamond, W. D., Chang, P. P., Folsom, A. R. & Chambless, L. E. Heart  
498 failure incidence and survival (from the Atherosclerosis Risk in Communities study). *Am. J.*  
499 *Cardiol.* **101**, 1016–1022 (2008).
- 500 10. Bui, A. L., Horwich, T. B. & Fonarow, G. C. Epidemiology and risk profile of heart failure.  
501 *Nat. Rev. Cardiol.* **8**, 30–41 (2011).
- 502 11. Roizen, M. F. Forecasting the Future of Cardiovascular Disease in the United States: A  
503 Policy Statement From the American Heart Association. *Yearbook of Anesthesiology and*  
504 *Pain Management* **2012**, 12–13 (2012).
- 505 12. Yancy, C. W., Jessup, M., Bozkurt, B. & Butler, J. 2013 ACCF/AHA guideline for the  
506 management of heart failure: a report of the American College of Cardiology  
507 Foundation/American Heart Association Task Force on .... *Journal of the* (2013).
- 508 13. Huang, H. *et al.* Accuracy of left ventricular ejection fraction by contemporary multiple  
509 gated acquisition scanning in patients with cancer: comparison with cardiovascular  
510 magnetic resonance. *Journal of Cardiovascular Magnetic Resonance* **19**, (2017).
- 511 14. Pellikka, P. A. *et al.* Variability in Ejection Fraction Measured By Echocardiography, Gated  
512 Single-Photon Emission Computed Tomography, and Cardiac Magnetic Resonance in  
513 Patients With Coronary Artery Disease and Left Ventricular Dysfunction. *JAMA Netw Open*  
514 **1**, e181456 (2018).
- 515 15. Malm, S., Frigstad, S., Sagberg, E., Larsson, H. & Skjaerpe, T. Accurate and reproducible  
516 measurement of left ventricular volume and ejection fraction by contrast echocardiography:  
517 a comparison with magnetic resonance imaging. *J. Am. Coll. Cardiol.* **44**, 1030–1035

- 518 (2004).
- 519 16. Cole, G. D. *et al.* Defining the real-world reproducibility of visual grading of left ventricular  
520 function and visual estimation of left ventricular ejection fraction: impact of image quality,  
521 experience and accreditation. *Int. J. Cardiovasc. Imaging* **31**, 1303–1314 (2015).
- 522 17. Koh, A. S. *et al.* A comprehensive population-based characterization of heart failure with  
523 mid-range ejection fraction. *Eur. J. Heart Fail.* **19**, 1624–1634 (2017).
- 524 18. Chioncel, O. *et al.* Epidemiology and one-year outcomes in patients with chronic heart  
525 failure and preserved, mid-range and reduced ejection fraction: an analysis of the ESC Heart  
526 Failure Long-Term Registry. *Eur. J. Heart Fail.* **19**, 1574–1585 (2017).
- 527 19. Shah, K. S. *et al.* Heart Failure With Preserved, Borderline, and Reduced Ejection Fraction:  
528 5-Year Outcomes. *J. Am. Coll. Cardiol.* **70**, 2476–2486 (2017).
- 529 20. Papolos, A., Narula, J., Bavishi, C., Chaudhry, F. A. & Sengupta, P. P. US hospital use of  
530 echocardiography: insights from the nationwide inpatient sample. *J. Am. Coll. Cardiol.* **67**,  
531 502–511 (2016).
- 532 21. ACCF/ASE/AHA/ASNC/HFSA/HRS/SCAI/SCCM/SCCT/SCMR 2011 Appropriate Use  
533 Criteria for Echocardiography. *J. Am. Soc. Echocardiogr.* **24**, 229–267 (2011).
- 534 22. Zhang, J. *et al.* Fully automated echocardiogram interpretation in clinical practice:  
535 feasibility and diagnostic accuracy. *Circulation* **138**, 1623–1635 (2018).
- 536 23. Madani, A., Arnaout, R., Mofrad, M. & Arnaout, R. Fast and accurate view classification of  
537 echocardiograms using deep learning. *NPJ Digit Med* **1**, (2018).
- 538 24. Ghorbani, A., Ouyang, D., Abid, A., He, B. & Chen, J. H. Deep Learning Interpretation of  
539 Echocardiograms. *bioRxiv* (2019).
- 540 25. Ardila, D. *et al.* Author Correction: End-to-end lung cancer screening with three-

- 541 dimensional deep learning on low-dose chest computed tomography. *Nat. Med.* **25**, 1319  
542 (2019).
- 543 26. Gulshan, V. *et al.* Development and Validation of a Deep Learning Algorithm for Detection  
544 of Diabetic Retinopathy in Retinal Fundus Photographs. *JAMA* **316**, 2402–2410 (2016).
- 545 27. Poplin, R. *et al.* Prediction of cardiovascular risk factors from retinal fundus photographs  
546 via deep learning. *Nat Biomed Eng* **2**, 158–164 (2018).
- 547 28. Esteva, A. *et al.* Dermatologist-level classification of skin cancer with deep neural  
548 networks. *Nature* **546**, 686 (2017).
- 549 29. Coudray, N. *et al.* Classification and mutation prediction from non–small cell lung cancer  
550 histopathology images using deep learning. *Nat. Med.* **24**, 1559–1567 (2018).
- 551 30. Chen, L.-C., Papandreou, G., Schroff, F. & Adam, H. Rethinking Atrous Convolution for  
552 Semantic Image Segmentation. *arXiv [cs.CV]* (2017).
- 553 31. Tran, D. *et al.* A Closer Look at Spatiotemporal Convolutions for Action Recognition. *2018*  
554 *IEEE/CVF Conference on Computer Vision and Pattern Recognition* (2018).  
555 doi:10.1109/cvpr.2018.00675
- 556 32. Tran, D., Bourdev, L., Fergus, R., Torresani, L. & Paluri, M. Learning spatiotemporal  
557 features with 3d convolutional networks. in *Proceedings of the IEEE international*  
558 *conference on computer vision* 4489–4497 (2015).
- 559 33. Kay, W. *et al.* The Kinetics Human Action Video Dataset. *arXiv [cs.CV]* (2017).