

1 **A Novel Digital Twin Strategy to Examine the Implications of Randomized Clinical Trials**  
2 **for Real-World Populations**

3  
4 Phyllis M. Thangaraj<sup>1\*</sup>, Sumukh Vasisht Shankar<sup>1\*</sup>, Sicong Huang<sup>3</sup>, Girish N. Nadkarni<sup>4,5</sup>, Bobak  
5 J. Mortazavi<sup>3</sup>, Evangelos K. Oikonomou<sup>1</sup>, and Rohan Khera<sup>1,2,6,7</sup>

6  
7 <sup>1</sup>Section of Cardiovascular Medicine, Department of Internal Medicine, Yale School of  
8 Medicine, New Haven, CT, USA

9 <sup>2</sup>Section of Health Informatics, Department of Biostatistics, Yale School of Public Health,  
10 New Haven, CT

11 <sup>3</sup>Department of Computer Science and Engineering, Texas A&M University, College Station,  
12 TX

13 <sup>4</sup>The Charles Bronfman Institute for Personalized Medicine, Icahn School of Medicine at  
14 Mount Sinai, New York, NY, USA

15 <sup>5</sup>The Division of Data Driven and Digital Medicine, Department of Medicine, Icahn School of  
16 Medicine at Mount Sinai, New York, NY, USA

17 <sup>6</sup>Section of Biomedical Informatics and Data Science, Yale School of Medicine, New Haven,  
18 CT

19 <sup>7</sup>Center for Outcomes Research and Evaluation, Yale-New Haven Hospital, New Haven, CT,  
20 USA

21 \*Contributed Equally

22  
23 Manuscript Type: Article

24 Words:

25 Abstract: 279

26 Text: 3824

27

28

29 **\*Address for correspondence:**

30 Rohan Khera, MD, MS

31 195 Church St, 6<sup>th</sup> Floor, New Haven, CT 06510

32 203-764-5885; rohan.khera@yale.edu; @rohan\_khera

33

34 **ABSTRACT**

35

36 Randomized clinical trials (RCTs) are essential to guide medical practice; however, their  
37 generalizability to a given population is often uncertain. We developed a statistically informed  
38 Generative Adversarial Network (GAN) model, RCT-Twin-GAN, that leverages relationships  
39 between covariates and outcomes and generates a digital twin of an RCT (RCT-Twin)  
40 conditioned on covariate distributions from a second patient population. We used RCT-Twin-  
41 GAN to reproduce treatment effect outcomes of the Systolic Blood Pressure Intervention Trial  
42 (SPRINT) and the Action to Control Cardiovascular Risk in Diabetes (ACCORD) Blood  
43 Pressure Trial, which tested the same intervention but found different treatment effects. To  
44 demonstrate treatment effect estimates of each RCT conditioned on the other RCT's patient  
45 population, we evaluated the cardiovascular event-free survival of SPRINT digital twins  
46 conditioned on the ACCORD cohort and vice versa (ACCORD twins conditioned on SPRINT).  
47 The conditioned digital twins were balanced across intervention and control arms (mean absolute  
48 standardized mean difference (MASMD) of covariates between treatment arms 0.019 (SD  
49 0.018), and the conditioned covariates of the SPRINT-Twin on ACCORD were more similar to  
50 ACCORD than SPRINT (MASMD 0.0082 SD 0.016 vs. 0.46 SD 0.20). Notably, across  
51 iterations, SPRINT conditioned ACCORD-Twin datasets reproduced the overall non-significant  
52 effect size seen in ACCORD (5-year cardiovascular outcome hazard ratio (95% confidence  
53 interval) of 0.88 (0.73-1.06) in ACCORD vs. median 0.87 (0.68-1.13) in the SPRINT  
54 conditioned ACCORD-Twin), while the ACCORD conditioned SPRINT-Twins reproduced the  
55 significant effect size seen in SPRINT (0.75 (0.64-0.89) vs. median 0.79 (0.72-0.86)) in the  
56 ACCORD conditioned SPRINT-Twin). Finally, we demonstrate the translation of this approach  
57 to real-world populations by conditioning the trials on an electronic health record population.  
58 Therefore, RCT-Twin-GAN simulates the direct translation of RCT-derived treatment effects  
59 across various patient populations.

## 60 INTRODUCTION

61 Randomized clinical trials (RCTs) generate evidence that defines optimal clinical practices, but  
62 their generalizability to real-world patient populations is often challenging to quantify.<sup>1,2</sup> This is  
63 a concern because RCTs often have underrepresentation from several demographic and clinical  
64 subpopulations<sup>3-7</sup> and varying treatment effects among individuals with certain characteristics.<sup>8-</sup>  
65 <sup>10</sup> These considerations are critical to translating information from RCTs to real-world patient  
66 populations,<sup>11,12</sup> but no strategies exist to evaluate how they may affect the applicability to  
67 patients in these settings.

68 Variation across RCTs testing similar interventions with discrepant treatment effects is a  
69 key issue for the generalizability of interventions tested in RCTs.<sup>13-19</sup> For example, the Systolic  
70 Blood Pressure Intervention Trial (SPRINT) was a treatment intervention RCT that showed  
71 improved cardiovascular outcomes with intensive blood pressure control.<sup>13</sup> In contrast, the  
72 Action to Control Cardiovascular Risk in Diabetes Blood Pressure (ACCORD) trial did not find  
73 improved cardiovascular outcomes with the same intervention.<sup>14</sup> Among the explanations posited  
74 for these discrepant findings include differences in population composition and event rates.<sup>20-23</sup>  
75 Despite experimental evidence from two trials, there is no quantitative strategy to evaluate these  
76 assertions explicitly. Therefore, while it is critical to evaluate whether the effects observed in an  
77 RCT population generalize to a second population – either a planned second RCT or a general  
78 population of patients with the condition – the challenge remains to examine these effects in the  
79 context of the complex differences across multiple population characteristics.

80 Digital twins of RCTs introduce a strategy to create a synthetic representation of a  
81 clinical trial updated by attributes of a second population. Specifically, trial-level digital twin  
82 synthesis through deep generative models such as Generative Adversarial Networks (GANs) can

83 integrate multiple covariates from a patient cohort by constructing a digital twin with covariate  
84 values sampled from the second cohort while retaining relationships and correlations between  
85 variables within the original RCT. While GANs have been utilized to estimate individual  
86 treatment effects, their potential for evidence translation across patient populations has not been  
87 explored.<sup>24–27</sup> Conditional GANs (CGAN) enable the generation of synthetic datasets that  
88 condition a model with covariates from a second population distribution.<sup>28,29</sup> We hypothesize that  
89 applying this model to an RCT conditioned on a second population will estimate the treatment  
90 effects of the original RCT in the new patient population.

91 We present RCT-Twin-GAN, a generative framework that combines clinical knowledge  
92 and the statistically informed architecture to create a digital twin of an RCT conditioned on the  
93 characteristics of a second patient population to assess for the generalizability of the treatment  
94 effect (Figure 1, Figure 2). To demonstrate the ability of the digital twin to replicate treatment  
95 effects in the conditioning population, we first compared two RCTs, SPRINT and ACCORD,  
96 with similar interventions but disparate treatment effects on cardiovascular outcomes. We created  
97 a digital twin of each of the 2 RCTs conditioned on covariate distributions of the other and  
98 evaluated whether the RCT-Twins reproduced the treatment effect of the conditioning cohort.  
99 Finally, we describe the cardiovascular outcomes of SPRINT and ACCORD digital twins  
100 conditioned on characteristics of patients in the electronic health record (EHR), introducing the  
101 role of RCT-Twins in estimating RCT treatment effects in real-world populations.

102

## 103 **RESULTS**

### 104 **Study Populations**

105 The study developed digital twins of two RCTs. The first RCT, SPRINT, was a treatment  
106 intervention study to test whether intensive blood pressure control (goal systolic blood pressure  
107 less than 120 mmHg) versus standard care (goal systolic blood pressure less than 140 mmHg)  
108 reduced major cardiovascular events. The trial consisted of 9361 participants (median age 67 (61  
109 to 76 (25-75% IQR, and 3332 (36%) women). The patients in SPRINT were followed for a  
110 median of 3.26 years for the first occurrence of any of the primary composite outcome of  
111 myocardial infarction, acute coronary syndrome, stroke, heart failure, or death from  
112 cardiovascular cause.

113 Our study built a SPRINT digital twin with a population representation of another RCT  
114 with the same intervention, the ACCORD trial, a double factorial RCT of participants with type  
115 2 diabetes mellitus and cardiovascular disease. We specifically leveraged the blood pressure  
116 management component of the ACCORD trial, wherein half of the participants were randomized  
117 to intensive versus standard care blood pressure control, with the same treatment goals as those  
118 in the SPRINT trial. ACCORD consisted of 4733 participants (median age 62, IQR, 58-67, and  
119 2258 [48%] women). ACCORD median follow-up time was 4.7 years for the primary composite  
120 outcome of myocardial infarction, stroke, or death from cardiovascular cause.

121 We also incorporated two cohorts from the Yale New Haven Hospital Health System  
122 Electronic Health Record (EHR), a large healthcare system including several hospitals with  
123 diverse racial and socioeconomic demographics across Connecticut and Rhode Island. Two sets  
124 of patients with hypertension, one without (N=22,132) and the other with diabetes (N=8,840)  
125 were identified to broadly represent populations included in SPRINT and ACCORD,  
126 respectively, to estimate the treatment effects found in the two RCTs on corresponding real-  
127 world patient populations. The final cohorts included 3,130 patients in the SPRINT EHR cohort

128 and 2,731 patients in the ACCORD EHR cohort. The SPRINT EHR cohort had a median age of  
129 73 years (IQR, 61 to 84) and 2069 (52%) women), while the ACCORD EHR cohort had a  
130 median age of 71 (IQR, 61 to 80) and 2032 (51%) women).

131

### 132 **The Non-Conditioned SPRINT Digital Twin Cohort**

133 We created 10 SPRINT-Twins (the non-conditioned SPRINT twin), which had a median age of  
134 66 (IQR, 60 to 75) and 1516-1704 (32-38%) women (Table S1, S2). The SPRINT-Twin  
135 reproduced the distributions of the original variables (covariates, outcome, and time to outcome)  
136 in SPRINT as evidenced by an absolute standardized mean difference (ASMD) of less than 0.1  
137 for each variable and a mean absolute standardized difference (MASMD) of 0.020 (SD 0.015)  
138 between the SPRINT Control (C) Arm and SPRINT-Twins C Arm and 0.021 (SD 0.014) between  
139 the SPRINT Intervention (I) Arm and SPRINT-Twins I Arm. In addition, all variables were  
140 balanced between the I and C arms in the SPRINT-Twin, as evidenced by an ASMD of less than  
141 0.1 for each variable and a MASMD of 0.011 (SD 0.016) between treatment arms across all  
142 variables. This was similar to the MASMD between treatment arms of SPRINT, 0.021 (SD  
143 0.018) and below the threshold where distributions are considered substantially dissimilar. The  
144 correlations between variables were also preserved as evidenced by 88.4% concordance between  
145 the Spearman correlations calculated between SPRINT's variables and those calculated between  
146 the SPRINT twin's variables (Table S3, Figure S1).

147

### 148 **The Conditioned SPRINT<sub>ACCORD</sub> and ACCORD<sub>SPRINT</sub> Digital Twin Cohorts**

149 We then generated 10 SPRINT<sub>ACCORD</sub> Twins, which were SPRINT twins conditioned with values  
150 from the ACCORD cohort for 10 covariates, and 10 ACCORD<sub>SPRINT</sub> Twins, which were

151 ACCORD twins conditioned with values from the SPRINT cohort for 10 covariates. The  
152 SPRINT<sub>ACCORD</sub> Twins had a median age of 62 years (IQR 58 to 68), 1106-1178 (46-49%) women  
153 (Tables S4, S5) with mean 2345 (SD 25.9) or 49.5% in the C arm, 2388 (SD 24) or 50.5% in the  
154 I arm. ACCORD<sub>SPRINT</sub> Twins had a median age of 67 years (IQR 61 to 76), 1545-1677 (32-34%)  
155 women with mean 4759 (SD 66) or 50.8% in the C arm and 4603 (SD 62) or 49.2% in the I arm  
156 (Tables S6, S7). Across all treatment arm covariate distributions of the SPRINT<sub>ACCORD</sub> Twins and  
157 ACCORD<sub>SPRINT</sub> Twins, there was little difference between the I and C arms, suggesting balanced  
158 treatment arms as evidenced by each covariate having an ASMD between treatment arms of less  
159 than 0.1 with an MASMD between treatment arms of 0.024 (SD 0.017) for SPRINT<sub>ACCORD</sub>  
160 Twins and 0.018 (SD 0.004) for ACCORD<sub>SPRINT</sub> Twins, respectively (Figure 3a).

161 Comparing datasets, across all the conditioned covariates, the ASMD between the  
162 SPRINT<sub>ACCORD</sub>-Twin and ACCORD were less than 0.1, with a MASMD of 0.008 (SD 0.016),  
163 and similarly, the ASMDs between ACCORD<sub>SPRINT</sub> Twin and SPRINT were less than 0.1 with a  
164 MASMD 0.023 (SD 0.014) compared to an MASMD of 0.46 (SD 0.20) for the same covariates  
165 between SPRINT and ACCORD (Figure 3b). Out of the six non-conditioned covariates, white  
166 race, systolic blood pressure, smoker, and LDL cholesterol level had ASMDs less than 0.1  
167 between ACCORD vs. SPRINT<sub>ACCORD</sub> Twin while systolic blood pressure, smoker, and angina  
168 had ASMDs less than 0.1 between SPRINT and ACCORD<sub>SPRINT</sub> Twin (Figure 3b). Conversely,  
169 when conditioning on the opposite RCT and comparing datasets (ie. ACCORD vs  
170 ACCORD<sub>SPRINT</sub> Twin and SPRINT vs SPRINT<sub>ACCORD</sub> Twin), the ASMD resembles ACCORD vs  
171 SPRINT for the conditioned covariates (Figure S2). Similar to the non-conditioned twins, the  
172 correlations between variables were also preserved, as evidenced by the 85.6% concordance of  
173 the Spearman correlations between the ACCORD variables and the correlations between the

174 SPRINT<sub>ACCORD</sub> Twin variables, the 78.4% concordance of correlations between the SPRINT  
175 variables and the correlations between the SPRINT<sub>ACCORD</sub> Twin variables, the 84.5%  
176 concordance of the correlations between the ACCORD variables and the ACCORD<sub>SPRINT</sub> twin  
177 variables, and the 85.6% concordance of the correlations between the SPRINT variables and the  
178 correlations between the ACCORD<sub>SPRINT</sub> Twin variables (Table S3, Figure S1).

179

### 180 **Digital Twin Similarity Evaluation**

181 Given the generated nature of the complementary covariates, each row of the conditioned twins  
182 does not perfectly match the original cohort patients since it is updated with the conditioning  
183 cohort data, so covariate distribution level assessments were conducted. When training a  
184 multivariate logistic regression classifier to distinguish between RCT and twin data, in which an  
185 accuracy of 0.5 is considered random chance, we found the median accuracy of the model to  
186 correctly classify the data as real or fake to be 0.50 (IQR 0.49 to 0.51) for distinguishing  
187 SPRINT from SPRINT Twins, 0.50 (IQR 0.49 to 0.51) for distinguishing SPRINT from  
188 SPRINT<sub>ACCORD</sub> Twins, 0.50 (IQR 0.49 to 0.55) for distinguishing ACCORD from ACCORD  
189 Twins, and 0.50 (IQR 0.50 to 0.51) for distinguishing ACCORD from ACCORD<sub>SPRINT</sub> Twins,  
190 compatible with the SPRINT twins not being distinguishable from the original trial. When  
191 assessing differentiation capability for each covariate by training and testing single variate  
192 logistic regression models, the overall median accuracy was 0.50-0.51 across comparisons  
193 (Figure S3).

194

### 195 **Sensitivity Analyses**



196 In a sensitivity analysis generating the SPRINT<sub>ACCORD</sub>-Twin, we assessed the convergence of the  
197 models at various training sample sizes, batch sizes, and number of epochs of training. At a  
198 training size of 1% and 25% of the SPRINT data, 81% of the models converged, while 94%  
199 converged with 50% and 100% of the data. All models were balanced across the training sample  
200 sizes. The proportion of models that reproduced the non-significant hazard ratio seen in  
201 ACCORD increased with the training sample and was in a majority of the samples with a sample  
202 size >10% of the trial (Table S8).

203

#### 204 **Comparison of RCT-Twin-GAN to other Synthesizer Models**

205 Our method consistently scored among the best in all statistical comparisons and correlations  
206 (Table S9, S10). It was superior to the other methods in machine learning efficacy, in which a  
207 gradient boosting classifier trained on generated digital twin values predicted original RCT  
208 values (Table S11).

209

#### 210 **Estimating the Primary Composite Outcome in the Non-Conditioned Twins**

211 We confirmed the differences in the reported primary composite outcomes in the SPRINT and  
212 ACCORD trials, which included a significant reduction in cardiovascular events in SPRINT's  
213 intervention arm compared with control (hazard ratio 0.75 [0.64-0.89 95% CI, p<0.001]) without  
214 a significant reduction in a similar primary composite outcome in ACCORD (hazard ratio 0.88  
215 [0.73 to 1.06 95% CI, p=0.20]). In the SPRINT-Twin without conditioning, the median hazard  
216 ratio across 10 generated SPRINT-Twin datasets was 0.73 (CI 0.61-0.87), with the 10  
217 replications performed to ensure the reproducibility of the findings. This was comparable to the  
218 HR of 0.75 in the SPRINT trial.<sup>13</sup> Similarly, the ACCORD-Twin without conditioning replicated

219 the primary results of the ACCORD trial, with a median HR of 0.89 (CI 0.79-1.0) comparable to  
220 the HR of 0.88 of the ACCORD trial.<sup>14</sup>

221

### 222 **Estimating the Primary Composite Outcome in the Conditioned Twins**

223 We then demonstrated the ability of RCT-Twins to replicate the known treatment effects of a  
224 second population with the  $\text{SPRINT}_{\text{ACCORD-Twin}}$  – the SPRINT-Twin that was conditioned on  
225 ACCORD. We found the median hazard ratio of 10  $\text{SPRINT}_{\text{ACCORD-Twin}}$  datasets was 0.87 (CI  
226 0.68-1.13), this time comparable to the HR of 0.88 of the ACCORD trial (Figure 4a). In contrast,  
227 in 10 replicated digital twins of the ACCORD cohort conditioned on covariate distributions in  
228 SPRINT ( $\text{ACCORD}_{\text{SPRINT-Twin}}$ ), reproduced the significant effect size seen in SPRINT (HR  
229 0.75) with a median hazard ratio of 0.79 (CI 0.72-0.86) (Figure 4b).

230

### 231 **Estimating the Treatment Effect of SPRINT and ACCORD in the EHR**

232 In a descriptive substudy, we demonstrated the ability to estimate SPRINT and ACCORD  
233 primary composite outcomes in patient populations reflecting a large US health system, YNHHS.  
234 The same 10 covariates used to build conditioned SPRINT and ACCORD twins were  
235 computably extracted from the YNHHS EHR by clinician experts to define covariates in the  
236 corresponding EHR cohorts and build digital twins of SPRINT and ACCORD conditioned on  
237 corresponding EHR cohorts (Tables S12-S15). In the digital twin of SPRINT conditioned on the  
238 corresponding EHR cohort ( $\text{SPRINT}_{\text{EHR-Twin}}$ ), we confirmed the replication of RCT features,  
239 including covariate balance across treatment arms (MASMD 0.03 (SD 0.03), (Figure S4). In this  
240  $\text{SPRINT}_{\text{EHR-Twin}}$  the median primary composite outcome HR was 0.84 (95% CI, 0.64-1.09)

241 across the 10 replications. Similarly, the ACCORD<sub>EHR</sub>-Twin replicated both RCT features and  
242 EHR covariate distributions, with a median primary composite outcome HR of 0.94 (CI 0.8-1.1).

243

## 244 **DISCUSSION**

245

246 We present RCT-Twin-GAN, a deep generative model that utilizes clinical knowledge of  
247 covariate relationships to synthesize a digital twin of an RCT with selected covariate  
248 distributions from a second population, such as another RCT cohort or a general patient  
249 population reflected in an EHR. RCT-Twin-GAN created digital twins that replicate the  
250 fundamental feature of RCTs, i.e., balanced covariates across treatment arms, but also reflected  
251 the covariate distributions of this second population's distribution. In addition, the RCT-Twin-  
252 GAN digital twin cohorts were indistinguishable from the SPRINT RCT cohorts, reproduced  
253 RCT covariate correlations, and outperformed other model architectures. Moreover, in a positive  
254 control experiment within a two-RCT system where treatment effects were known from well-  
255 conducted experiments but were discordant across the RCTs, the RCT-Twins conditioned on  
256 covariates from the opposing RCT replicated the results observed in the other RCT,  
257 demonstrating the value of examining the effect of population characteristics on study outcomes.  
258 We also demonstrate that the approach is flexible to these characteristics drawn from any  
259 population, thereby enabling a quantitative evaluation of an RCT's potential treatment effects in  
260 populations that differed from those included in the trial.

261 Our work has built upon the established need to quantify generalizability of RCTs to new  
262 populations.<sup>32</sup> Prior methods, such as standardization of event rates, allow adjustment by single  
263 variables, which groups patients together by singular stratification.<sup>33</sup> Others have used distance  
264 metrics and decision tree machine learning techniques to represent the complex interplay of

265 covariates and characterize the heterogeneity of treatment effect.<sup>8–10,23,34–36</sup> Prior generative  
266 methods have used statistical machine learning to build digital twins of control patients in  
267 neurological clinical trials and observational studies with accurate reproduction of patient  
268 trajectory at the individual level.<sup>35–37</sup> Our method complements these by building trial-level  
269 digital twins of a conditioning cohort that draw from the multiple covariate distributions and  
270 outcomes of an RCT population to generate equivalent covariates in the conditioning population  
271 and estimate population-level treatment effects from the RCT intervention. CTAB-GAN+ has  
272 been used to build an RCT control patient population, but our studies have demonstrated  
273 superiority with DATGAN in reproducing trial baseline characteristics.<sup>38</sup> In addition, compared  
274 to other GAN conditioning methods, our architecture is the only method that can condition on  
275 multiple continuous and categorical variables, allowing for multi-variate correlations of the  
276 conditioning cohort to be preserved. Statistical methods to assess heterogeneous treatment effects  
277 across populations have generally focused on equalizing baseline characteristics between  
278 populations using propensity score matching, but this scores one variable at a time, thereby  
279 ignoring multi-variable differences across patients, and does not consider effect modifiers.<sup>39</sup>

280 We incorporate the distributions of multiple mutual pre-randomization covariates  
281 available across datasets to ensure representation across multivariate axes. In addition, we utilize  
282 clinician expertise to identify connections between covariates and build digital twins modeling  
283 the complex interplay of effect modifiers and outcomes. The result is a data-driven generated  
284 outcome of the conditioning cohort based on the correlations between multiple covariates and  
285 within overlapping covariate distributions between the two patient populations. In Table 1, we  
286 discuss the minimum requirements to estimate treatment effects across two populations,  
287 including cohort requirements, randomization, intervention and outcome, and sample size.

288           A unique feature of our model incorporates both rigorous statistical methods and clinical  
289 knowledge to build digital twins of RCTs with representative covariate balance and effect  
290 modifier information. The DAG structure weights clinically relevant relationships between  
291 covariates and outcomes and removes spurious correlations that would otherwise be included in  
292 the GAN. Of note, the choice of the covariates was governed by primary analysis focused on  
293 shared covariates between SPRINT and ACCORD. In real-world translations, a different  
294 covariate set shared between a development and target population can be selected. In addition,  
295 the ability of the ciDATGAN architecture to condition on multiple continuous and categorical  
296 values is unique compared to competing architectures. Our ability to reproduce treatment effect  
297 estimates from the conditioning cohort by sampling its covariate distributions relies on the  
298 inference of important correlations between covariates during GAN training and digital twin  
299 generation. Although prior digital twin studies have focused on supplementing RCTs with  
300 synthetic patients for controls<sup>35-38</sup> and reproducing progression within the same cohort, our study  
301 builds upon these by estimating the treatment effect across different patient populations.  
302 Measuring the hazard ratios of treatment effect outcomes as an evaluation metric provided  
303 valuable insights into the fidelity of the synthetic dataset in simulating clinical trial outcomes and  
304 treatment responses.

305           Methodologically, ACCORD represents a second RCT that experimentally tested the  
306 same intervention as SPRINT but in a different population. This is essential as the effect  
307 estimates in a conditioned twin otherwise have no gold standard comparison. We demonstrate  
308 that conditioning generates effect estimates replicated in the trial that experimentally tested the  
309 intervention but arrived at a different conclusion, suggesting the validity of the produced  
310 estimates. This is the key methodological outcome of our experiments.

311           Moreover, there are direct clinical implications for both hypertension management and  
312 evidence generation via clinical trials. Our work demonstrates that the observed effect estimate  
313 differences in SPRINT and ACCORD emerge because of the nature of populations enrolled in  
314 these trials, and not because of diabetes status. There has been a lack of clarity about whether  
315 these effects suggested some consideration about blood pressure and its effects on diabetes. But  
316 we demonstrate that even in a trial like SPRINT, if the enrolled population had key features that  
317 resemble those seen in ACCORD, the trial could have produced a potentially null result.  
318 Similarly, had ACCORD enrolled patients that resembled SPRINT – based on features other than  
319 diabetes, it could have been positive. While these observations represent data experiments, this  
320 was recently observed in the ESPRIT trial, where patients with diabetes benefitted from intensive  
321 blood pressure lowering.<sup>40</sup>

322           This has implications for the interpretation of clinical trials as well. We acknowledge that  
323 our work is a proof-of-concept, but we demonstrate that trials can be evaluated in populations  
324 that differ from those enrolled on key features to address whether embedded heterogeneous  
325 treatment effects and differences in these covariates affect how these results should be  
326 interpreted. Moreover, these experiments can guide the need for populations ideally chosen for  
327 additional trials. Health systems could determine the likely treatment effect of an intervention in  
328 their patient population to better contextualize their patient outcomes with the intervention by  
329 developing population-wide digital twins. This effort to use general real-world evidence to  
330 establish the efficacy of interventions has major regulatory support from agencies such as the US  
331 Food and Drug Administration.<sup>41</sup>

332           There are limitations to consider. First, RCT-Twin-GAN uses a select set of variables to  
333 build the digital twin. We chose a smaller set of covariates to maximize efficiency and showed

334 that even with this small number of representative variables, we can build a digital twin that  
335 successfully replicates treatment effect estimates. Second, our model relies on outside input for  
336 identifying correlations between covariates, but we believe this can be considered a strength that  
337 clinical expertise can be imbued into the model to reduce the weight of spurious correlations  
338 inherent in data. Third, this is a post-hoc analysis of RCTs, but we show the ability of digital  
339 twins to mirror covariate characteristics and treatment effects found in SPRINT and ACCORD.  
340 Fourth, we only applied RCT-Twin-GAN to the SPRINT - ACCORD pair because it was the only  
341 paired trial testing the same intervention with different results available through a public domain,  
342 the National Heart, Lung, and Blood Institute Biologic Specimen and Data Repository  
343 Information Coordinating Center (BioLINCC). As the data are publicly available, further  
344 research can build upon this example, and we further anticipate applying our model to other  
345 examples.

346 Fifth, we could not study glycemic effects on the intervention because SPRINT did not  
347 control hyperglycemia or include diabetic patients as seen in ACCORD. Despite this, we  
348 demonstrate a positive treatment effect aligned with SPRINT among the ACCORD patients, who  
349 had a full range of glycemic management differences as part of the original ACCORD trial, and  
350 so do not find evidence to suggest glycemic management differences produced the null results  
351 observed in ACCORD. Sixth, GANs are known to have challenges with achieving successful  
352 convergence between the discriminator and generator, so we have adapted the most successful  
353 advances in GAN development to optimize convergence. We chose an architecture that reduces  
354 spurious correlations by introducing a DAG to define the correlation structure between variables,  
355 stabilized training with the most appropriate learning rate and batch normalization layers, and  
356 incorporated a loss function that eliminated the risk of vanishing gradients to ensure optimal

357 model performance. Seventh, modeling real-world patients in the EHR can be challenging since  
358 the data represents a snapshot of patients who seek care, but we choose patients from a diverse  
359 tertiary care system to maximize the breadth of the general population identified. In addition, the  
360 EHR covariates had to be operationally defined by experts to be analogous to the criteria used in  
361 RCT, but this is a descriptive study that shows different covariate distributions can be modeled.  
362 Finally, the true effect estimates in the EHR populations are unknown, and those estimated by  
363 RCT-Twin-GAN should not inform care but rather give an idea of discordance or concordance  
364 with the original RCT population.

365 We have introduced a new application of GANs to build synthetic cohorts by creating an  
366 RCT digital twin reflective of different patient populations, including similar RCTs and real-  
367 world patients found in the EHR. Our study demonstrates a way to evaluate the generalizability  
368 of an RCT to the general population by embedding covariate distributions that are more  
369 representative of real-world populations. This amplifies the effects for those more frequently  
370 seen in clinical practice. Overall, our model contributes significantly to the evidence supporting  
371 the development of an RCT digital twin that more consistently mirrors real-world populations,  
372 thereby enhancing inference for real-world patients.

373

## 374 **METHODS**

375

### 376 **Data Source and Patient Populations**

#### 377 SPRINT and ACCORD Cohorts

378 From 2010-2013, at 102 clinical sites across the United States, participants were recruited for the  
379 SPRINT RCT who were at least 50 years old, had a systolic blood pressure between 130 and 180  
380 mm Hg, and had increased cardiovascular event risk, including cardiovascular disease with the



381 exception of stroke, chronic kidney disease, Framingham 10 year cardiovascular risk score of  
382 15% or greater, and advanced age over 75. Patients with prior stroke, diabetes mellitus, and a  
383 recent heart failure exacerbation had been excluded from the study.

384 From 2001 to 2005, at 77 clinical sites across the United States and Canada, participants  
385 were recruited for the ACCORD RCT who had type 2 diabetes mellitus, a glycated hemoglobin  
386 level of 7.5% or greater, and either age 40 or older with cardiovascular disease or age 55 or older  
387 with risk factors for cardiovascular disease and anatomical evidence of longstanding  
388 hypertension or diabetes such as albuminuria or left ventricular hypertrophy. Patients with a BMI  
389 over 45, a creatinine over 1.5 mg/dL, or serious illness were excluded.

#### 390 EHR cohorts

391 The two EHR cohorts were extracted from patients within the Yale New Haven Health System  
392 (YNHHS) from 2013 to 2023. The study was reviewed by the Yale Institutional Review Board  
393 and deemed exempt as it uses retrospective data. We sampled 100,000 adult patients and then  
394 filtered the cohort to those with an ICD-10-CDM code for hypertension (Table S16). Out of these  
395 patients, we filtered for patients with an ICD-10-CDM code for type 2 diabetes mellitus (Table  
396 S16). Patients with both hypertension and type 2 diabetes mellitus billing codes were considered  
397 for the ACCORD EHR cohort. The remaining hypertension patients who did not have type 2  
398 diabetes mellitus billing codes were considered for the SPRINT EHR cohort. We excluded  
399 patients who did not have values for continuous covariates and patients above the age of 110. We  
400 then sampled 4000 patients each for the ACCORD EHR and SPRINT EHR cohorts with values  
401 for all conditioned covariates. We further excluded patients who had continuous values out of  
402 range of the training cohort of SPRINT or ACCORD (Table S17).

403

## 404 **Development of RCT Digital Twins Conditioned on a Second Patient Population**

405 We adapted CGAN models to create digital twin datasets of an RCT conditioned on covariate  
406 distributions from a second patient population. We first built a SPRINT digital twin (SPRINT-  
407 twin) trained on the SPRINT cohort without a second conditioning cohort. We then built a  
408 SPRINT digital twin conditioned on the ACCORD participant population (SPRINT<sub>ACCORD</sub>-Twin)  
409 with the intention of reproducing the ACCORD primary composite outcome in a SPRINT digital  
410 twin (Figure 1). To implement this, we applied the Conditional inputs for Direct Acyclic Tabular  
411 Generative Adversarial Networks (CiDATGANs), a conditional tabular GAN that uses a directed  
412 acyclic graph (DAG) to assign relationships between pre-randomized covariates.<sup>29,42</sup> The DAG  
413 ensures clinically relevant connections are introduced between covariates and prevents the  
414 weighting of spurious correlations between covariates. To condition the digital twins on the other  
415 RCT population, we mapped 33 equivalent covariates between SPRINT and ACCORD (Table  
416 S18).

417

## 418 **Covariate Extraction for SPRINT, ACCORD, and the EHR**

419 In order to condition the SPRINT digital twin (SPRINT-Twin) on equivalent ACCORD  
420 covariates (SPRINT<sub>ACCORD</sub>-Twin), we mapped 33 equivalent covariates between the two cohorts,  
421 which included demographics such as age, gender, race, and ethnicity, conditions and social  
422 history, such as smoking history, family history of cardiovascular disease (CVD),  
423 hyperlipidemia, left ventricular hypertrophy (LVH) and prior myocardial infarction (MI),  
424 medications such as taking aspirin or statins, procedures such as coronary revascularization, and  
425 laboratory values and vital signs such as glomerular filtration rate (GFR), glucose, and systolic  
426 blood pressure (Table S18). We also included outcome, time to outcome, and treatment arm

427 assignment. We limited the maximum time to outcome to five years, censoring all subsequent  
428 outcomes.

429 To build the DAG, an expert clinician identified 16 representative variables of the 33  
430 mapped between SPRINT and ACCORD to represent all clinical areas such as demographics,  
431 conditions, medications, family history, symptoms, social history, procedures, vital signs, and  
432 laboratory and EKG measures, and also maintaining a balance of both categorical and continuous  
433 variables. All demographic variables were included since they are available for everyone in the  
434 EHR cohort. The variables included continuous covariates of age at randomization, GFR, heart  
435 rate, LDL cholesterol, and systolic blood pressure, and categorical covariates converted to a  
436 binary assignment of the presence (1) or absence (0) of angina, Black race, BMI, current smoker,  
437 family history of CVD, female sex, Hispanic ethnicity, LVH, previous MI, statin use, and White  
438 race (Table S18). Since BMI was considered a binary variable in ACCORD (above or below 32  
439 kg/m<sup>2</sup>), we used a similar definition in SPRINT.

440 Variables related to exclusion criteria of at least one of the cohorts were not included in  
441 the conditioning of the model or constructing the DAG because of the lack of overlap in the  
442 distribution of these covariate values between the SPRINT and ACCORD cohorts. These  
443 included glucose and diabetes mellitus. The DAG construction includes an iterative process of  
444 expert assessment of clinically relevant pairs and the causal direction within the pairs and  
445 calculation of correlations between unpaired variables (Table 2). The final DAG included 71  
446 connections (Figure 2, Table S19). The arrows' direction pointed from the independent covariate  
447 to the dependent covariate. No arrow pointed to the treatment arm covariate, labeled "Group",  
448 since this assignment was independent of all covariates. All covariates and the "Group" pointed  
449 to the "Outcome" and "Time to Outcome" covariates since all covariates and treatment arm

450 assignment were thought to influence the outcome (Figure 2, Table S19). We used the 10  
451 covariates with the largest absolute standardized mean difference between SPRINT and  
452 ACCORD as the conditioned covariates in order to condition from the covariate distributions  
453 most representative of the second cohort. The included binary and continuous covariates, in the  
454 order of increasing dissimilarity between cohorts, were black race, history of previous MI,  
455 female sex, statin use, LVH, BMI, heart rate, age at randomization, family history of CVD, and  
456 GFR.

457         Since we sought to condition on the EHR populations as well, we extracted the 10  
458 conditioned covariates established in the prior analysis from the EHR as well (Table S16). Only  
459 patients with a value for the demographics sex, race, and age (based on an available date of  
460 birth), the vital signs BMI and heart rate, and the laboratory test eGFR (or computable from  
461 serum creatinine), were included. The binary covariates of family history of CVD, LVH,  
462 previous MI, and statin use were considered not present (assigned 0) if they were not recorded in  
463 the patient's EHR, as is the norm for observational research studies in the EHR.<sup>43</sup> Age was  
464 calculated on October 1, 2023 (EHR query date), unless they were deceased, where we used the  
465 death date to define their last known age. We used this index date to consider most current  
466 clinical characteristics of the patient to estimate their treatment effect, the equivalent of the  
467 randomization to treatment arm date in the RCT.

468

### 469 **Design of the RCT-Twin-GAN Model**

470 RCT-Twin-GAN is a Generative Adversarial Network model, which is a deep learning model  
471 rooted in game theory that pits a generator, the neural network that creates synthetic data, against  
472 a discriminator, the neural network that distinguishes between the real data it is trained on and

473 the synthetic data created by the generator. The minimization of the discrimination between real  
474 and synthetic data allows for the GAN to make realistic digital twins of the cohort on which it is  
475 developed.<sup>44</sup> The neural networks are comprised of Long Short Term Memory (LSTM) cells,  
476 which are structured to retain information from prior inputs in addition to the current variable  
477 input.<sup>45</sup> GANs have been adapted to accurately synthesize tabular data such as EHR data.<sup>27,28,46</sup>  
478 GANs can also integrate data from a second patient population through conditioning the model  
479 on sample covariate values from the original cohort within the covariate distribution of the  
480 conditioning cohort, or the second patient population.<sup>28,29,46</sup> To avoid the well-documented  
481 challenges with consistently achieving convergence in GANs, our model utilizes Wasserstein  
482 loss to overcome training instability and prevent vanishing gradients.<sup>47</sup>

483 RCT-Twin-GAN is based on the architecture of CiDATGAN, which is an extension of  
484 DATGAN with an additional feature of conditioning covariates with distributions from a second  
485 population.<sup>29,42</sup> The DATGAN and CiDATGAN models employ a unique feature, allowing the  
486 generator to have the relationships between covariates and outcomes of the original training  
487 cohort to be explicitly encoded via a Directed Acyclic Graph (DAG). This prevents overfitting of  
488 the discriminator by defining the correlation structure between variables. It constrains the  
489 number of relevant associations, prioritizing key features for the model to learn. In contrast to  
490 other conditional GAN architectures that infer variable relationships solely by correlation, the  
491 addition of DAG to the training of the CiDATGAN generator incorporates directed relationships  
492 between pairs of variables to eliminate spurious correlations between variables. Continuous  
493 variables were winsorized based on the min-max values of the covariate in the training dataset  
494 (Table S17) to remove outlier values below the 2.5% and 97.5% percentiles, and categorical

495 variables were encoded into one-hot vectors and then fed into the discriminator as part of the  
496 input.

497         During the training phase, the generator combines Gaussian noise and attention vectors of  
498 the LSTM cells in the order of the DAG relationships and transforms the covariates from the  
499 original cohort using a fully connected layer in order to refine relationships and dependencies  
500 between the inputs. CiDATGAN creates a key modification to the DATGAN architecture in  
501 which the transformed conditional covariate inputs are also fed to the generator. Because of this  
502 modification, the DAG is also modified so that all conditional covariates are source nodes. The  
503 generator then synthesizes complementary values of the remaining covariates of the original  
504 dataset. The discriminator is then trained to differentiate between the original versus generated  
505 values of the remaining covariates from the original dataset. The discriminator is then trained to  
506 differentiate between the original versus generated values of the remaining covariates from the  
507 original dataset (Figure 1a).

508         During the sampling phase, the generator receives Gaussian noise, the modified DAG,  
509 and conditioned covariate values from the conditioning cohort and produces synthetic data  
510 without transformation in order to directly reflect the learned distribution from training while  
511 maintaining the integrity of the inputs from the conditioning cohort. Therefore, the final synthetic  
512 dataset incorporates the conditioning cohort inputs and generates complementary values for the  
513 remaining covariates missing from the conditioning cohort. The generator creates a cohort digital  
514 twin by producing one row of data at a time for each patient within the conditioning cohort  
515 (Figure 1b).

516         The CiDATGAN was then trained with the DAG and encoded dataset to generate the  
517 synthetic dataset. We performed a hyperparameter grid search with different batch sizes and

518 epochs to find the best parameters to generate synthetic data with similar outcomes to the  
519 original dataset (Table S20).

520

### 521 **Application of RCT-Twin-GAN in SPRINT, ACCORD, and the EHR**

522 We first used RCT-Twin-GAN to build a SPRINT-Twin, which created a DAG based on the  
523 SPRINT cohort, trained the DATGAN architecture on the SPRINT cohort, and ran the sampling  
524 phase of the DATGAN pipeline ten times, resulting in ten distinct synthetic twins. We then built  
525 a SPRINT<sub>ACCORD</sub>-Twin, which again created a DAG from and trained on the SPRINT cohort but  
526 was conditioned on the ACCORD cohort, utilizing the CiDATGAN architecture. This meant that  
527 the DAG was modified to remove connections going to the conditioning covariates, and in the  
528 sampling phase, Gaussian noise and the conditioned covariate distributions from the ACCORD  
529 cohort were inputs for the generator in order to create the final synthetic dataset. We sampled this  
530 process for 10 iterations to make 10 SPRINT<sub>ACCORD</sub>-Twin datasets. We repeated this process to  
531 create ACCORD<sub>SPRINT</sub>-Twin datasets by replacing ACCORD to be the training cohort and  
532 SPRINT to be the conditioning cohort.

533 We repeated this training and conditioning process using the EHR cohorts as well.

534 Specifically, we trained RCT-Twin-GAN on the SPRINT cohort, conditioned on the SPRINT-  
535 EHR cohort, and sampled 10 times to create SPRINT<sub>EHR</sub>-Twins. We similarly made  
536 ACCORD<sub>EHR</sub>-Twins with ACCORD and the ACCORD-EHR cohort.

537

### 538 **Analysis of Cohort Representation in Digital Twins**

539 To determine whether RCT-Twin-GAN created digital twins that are balanced by treatment arm,  
540 we calculated the mean absolute standardized mean difference (MASMD) of all covariates of the

541 digital twins stratified by treatment arm assignment. A value of less than 0.1 was considered  
542 adequate balance, consistent with convention when assessing the success of propensity score  
543 matching.<sup>48</sup> To assess the representation of the conditioning cohort in the synthetic digital twin,  
544 we also calculated the absolute standardized mean difference (ASMD) between SPRINT and  
545 ACCORD for each covariate, SPRINT-Twin and ACCORD, and SPRINT<sub>ACCORD</sub>-Twin and  
546 ACCORD. We also calculated the Spearman correlation between all variables for each RCT and  
547 digital twin, discretized the correlations into 7 bins from -1 to 1, and then calculated the  
548 proportion of covariate correlations from the RCT and twin data that were in the same bin  
549 (termed correlation accuracy) and mean absolute difference between RCT covariate correlation  
550 values and digital twin covariate correlation values as described in Li et al.<sup>27</sup>

551

### 552 **Evaluation of Digital Twin Similarity and Integrity**

553 We evaluated whether any row of the digital twin cohorts matched the RCT patients by a  
554 similarity score evaluation by Synthetic Data Vault.<sup>49</sup> To assess distinguishability of the RCT  
555 data from the conditioned twins, we trained and tested with a 70/30 split a multivariable logistic  
556 regression classifier to differentiate between the RCT and digital twin, which has been used to  
557 assess the integrity of prior digital twins.<sup>37</sup>

558

### 559 **Comparison of RCT-Twin-GAN to other synthesizer models**

560 We compared the DATGAN architecture used in RCT-Twin-GAN to 5 competing models  
561 including Conditional Tabular GAN (CTGAN)<sup>28</sup>, Conditional Tabular GAN+ (CTABGAN+)<sup>50</sup>,  
562 CopulaGAN<sup>51</sup>, GaussianCopula<sup>51</sup>, and Triplet-based Variational Autoencoder (TVAE)<sup>28</sup>. We  
563 utilized the non-conditioned DATGAN architecture because none of the comparator models have



564 the flexibility to condition multiple continuous and categorical variables like the CiDATGAN  
565 architecture. Specifically, we tested the mean absolute error,  $R^2$ , root mean squared error,  
566 standardized root mean squared error, and Pearson correlation between the distribution of unique  
567 values in the RCT compared to the digital twin. We also tested machine learning efficacy, in  
568 which a gradient boosting classifier is trained on twin data and evaluated on how well it generates  
569 real data.

570

### 571 **Estimation of Treatment Effect on Cardiovascular Outcomes in the Digital Twins**

572 In order to assess the ability of RCT-Twin-GAN to estimate RCT treatment effect outcomes in  
573 populations other than the original RCT, we calculated the hazard ratio of cardiovascular  
574 outcomes stratified by treatment arms in each of the digital twin cohorts using cox proportional  
575 hazard models. We utilized hazard ratios to evaluate the comparative risks of events over time  
576 between different treatment groups within the synthetic data. This analytical approach allowed us  
577 to gauge the effectiveness of the synthetic dataset in accurately representing the underlying  
578 dynamics of treatment effects and event occurrences observed in real-world scenarios.

579 We reported the median hazard ratio and 95% confidence intervals for the 10 SPRINT-  
580 Twin, SPRINT<sub>ACCORD</sub>-Twin, and ACCORD<sub>SPRINT</sub>-Twin digital twins. In order to demonstrate the  
581 ability to estimate treatment effect outcomes in a variety of cohorts, we calculated the hazard  
582 ratio and 95% confidence intervals of cardiovascular outcomes of the SPRINT<sub>EHR</sub>-Twins and  
583 ACCORD<sub>EHR</sub>-Twins as well.

584

### 585 **Statistical Analysis**

586 Categorical variables were summarized as numbers with percentages, and continuous variables  
587 were summarized as median with 25% and 75% interquartile ranges (IQR) or mean with  
588 standard deviation (SD). Covariate distributions were compared using ASMD and standard  
589 deviations and Spearman correlations between covariates and graphed as love plots comparing  
590 datasets and heatmaps of the correlations. Data was winsorized at 2.5% and 97.5% percentiles to  
591 remove outliers. Survival analysis was conducted using unadjusted cox proportional hazard  
592 models with p values calculated after 5 years and presented as Kaplan Meier survival curves.  
593 Hazard ratios across digital twins and SPRINT and ACCORD cohorts were presented as forest  
594 plots with 95% confidence interval error bars. Analyses were conducted using python 3.9, with  
595 packages specified in the supplement.

596

#### 597 **DATA AVAILABILITY**

598 The SPRINT and ACCORD cohorts are publicly available through the National Heart, Lung, and  
599 Blood Institute Biologic Specimen and Data Repository Information Coordinating Center  
600 (BioLINCC). The SPRINT dataset is available at <https://biolincc.nhlbi.nih.gov/studies/sprint/>  
601 and the ACCORD dataset is available at <https://biolincc.nhlbi.nih.gov/studies/accord/>. The Yale  
602 electronic health record cohorts are not available due to the use of patient data.

603

#### 604 **CODE AVAILABILITY**

605 The code for reproducing the treatment effect estimates, digital twins, and analysis figures will  
606 be available during peer review in an accompanying file, and the code will be made publicly  
607 available upon publication.

608

609

610 **REFERENCES**

- 611 1. Averitt AJ, Ryan PB, Weng C, Perotte A. A conceptual framework for external validity. *J*  
612 *Biomed Inform.* 2021;121:103870.
- 613 2. Rothwell PM. External validity of randomised controlled trials: “To whom do the results of  
614 this trial apply?” *Lancet.* 2005;365:82–93.
- 615 3. Filbey L, Zhu JW, D’Angelo F, et al. Improving representativeness in trials: a call to action  
616 from the Global Cardiovascular Clinical Trialists Forum. *Eur Heart J.* 2023;44:921–930.
- 617 4. Kennedy-Martin T, Curtis S, Faries D, Robinson S, Johnston J. A literature review on the  
618 representativeness of randomized controlled trial samples and implications for the external  
619 validity of trial results. *Trials.* 2015;16:495.
- 620 5. Ranganathan M, Bhopal R. Exclusion and inclusion of nonwhite ethnic minority groups in 72  
621 North American and European cardiovascular cohort studies. *PLoS Med.* 2006;3:e44.
- 622 6. Sardar MR, Badri M, Prince CT, Seltzer J, Kowey PR. Underrepresentation of women, elderly  
623 patients, and racial minorities in the randomized trials used for cardiovascular guidelines. *JAMA*  
624 *Intern Med.* 2014;174:1868–1870.
- 625 7. DeFilippis EM, Echols M, Adamson PB, et al. Improving Enrollment of Underrepresented  
626 Racial and Ethnic Populations in Heart Failure Trials: A Call to Action From the Heart Failure  
627 Collaboratory. *JAMA Cardiol.* 2022;7:540–548.

- 628 8. Oikonomou EK, Spatz ES, Suchard MA, Khera R. Individualising intensive systolic blood  
629 pressure reduction in hypertension using computational trial phenomaps and machine learning: a  
630 post-hoc analysis of randomised clinical trials. *Lancet Digit Health*. 2022;4:e796–e805.
- 631 9. Oikonomou EK, Suchard MA, McGuire DK, Khera R. Phenomapping-Derived Tool to  
632 Individualize the Effect of Canagliflozin on Cardiovascular Risk in Type 2 Diabetes. *Diabetes*  
633 *Care*. 2022;45:965–974.
- 634 10. Oikonomou EK, Van Dijk D, Parise H, et al. A phenomapping-derived tool to personalize the  
635 selection of anatomical vs. functional testing in evaluating chest pain (ASSIST). *Eur Heart J*.  
636 2021;42:2536–2548.
- 637 11. Patel HC, Hayward C, Dzung JN, et al. Assessing the Eligibility Criteria in Phase III  
638 Randomized Controlled Trials of Drug Therapy in Heart Failure With Preserved Ejection  
639 Fraction: The Critical Play-Off Between a “Pure” Patient Phenotype and the Generalizability of  
640 Trial Findings. *J Card Fail*. 2017;23:517–524.
- 641 12. Lim YMF, Molnar M, Vaartjes I, et al. Generalizability of randomized controlled trials in  
642 heart failure with reduced ejection fraction. *Eur Heart J Qual Care Clin Outcomes*. 2022;8:761–  
643 769.
- 644 13. SPRINT Research Group, Wright JT Jr, Williamson JD, et al. A Randomized Trial of  
645 Intensive versus Standard Blood-Pressure Control. *N Engl J Med*. 2015;373:2103–2116.
- 646 14. ACCORD Study Group, Cushman WC, Evans GW, et al. Effects of intensive blood-pressure  
647 control in type 2 diabetes mellitus. *N Engl J Med*. 2010;362:1575–1585.

- 648 15. Carson JL, Brooks MM, Hébert PC, et al. Restrictive or Liberal Transfusion Strategy in  
649 Myocardial Infarction and Anemia. *N Engl J Med*. 2023;389:2446–2456.
- 650 16. Ducrocq G, Gonzalez-Juanatey JR, Puymirat E, et al. Effect of a Restrictive vs Liberal Blood  
651 Transfusion Strategy on Major Cardiovascular Events Among Patients With Acute Myocardial  
652 Infarction and Anemia: The REALITY Randomized Clinical Trial. *JAMA*. 2021;325:552–560.
- 653 17. Joosten LPT, van Doorn S, van de Ven PM, et al. Safety of Switching from a Vitamin K  
654 Antagonist to a Non-Vitamin K Antagonist Oral Anticoagulant in Frail Older Patients with Atrial  
655 Fibrillation: Results of the FRAIL-AF Randomized Controlled Trial. *Circulation*. 2023.  
656 Published online August 27, 2023. <https://doi.org/10.1161/CIRCULATIONAHA.123.066485>.
- 657 18. Granger CB, Alexander JH, McMurray JJV, et al. Apixaban versus Warfarin in Patients with  
658 Atrial Fibrillation. *N Engl J Med*. 2011;365:981–992.
- 659 19. Jane-wit D, Horwitz RI, Concato J. Variation in results from randomized, controlled trials:  
660 stochastic or systematic? *J Clin Epidemiol*. 2010;63:56–63.
- 661 20. Krakoff LR. A tale of 3 trials: ACCORD, SPRINT, and SPS3. What happened? *Am J*  
662 *Hypertens*. 2016;29:1020–1023.
- 663 21. Chobanian AV. Hypertension in 2017-what is the right target? *JAMA*. 2017;317:579–580.
- 664 22. Huang C, Dhruva SS, Coppi AC, et al. Systolic blood pressure response in SPRINT (Systolic  
665 Blood Pressure Intervention Trial) and ACCORD (Action to Control Cardiovascular Risk in  
666 Diabetes): A possible explanation for discordant trial results. *J Am Heart Assoc*. 2017;6.

- 667 23. Laffin LJ, Besser SA, Alenghat FJ. A data-zone scoring system to assess the generalizability  
668 of clinical trial results to individual patients. *Eur J Prev Cardiol.* 2019;26:569–575.
- 669 24. Liu R, Rizzo S, Whipple S, et al. Evaluating eligibility criteria of oncology trials using real-  
670 world data and AI. *Nature.* 2021;592:629–633.
- 671 25. Ge Q, Huang X, Fang S, et al. Conditional Generative Adversarial Networks for  
672 Individualized Treatment Effect Estimation and Treatment Selection. *Front Genet.*  
673 2020;11:585804.
- 674 26. Yoon J, Jordon J, Van Der Schaar M. Ganite: Estimation of individualized treat- ment effects  
675 using generative adversarial nets. 2018. Accessed November 9, 2023.  
676 <https://openreview.net/pdf?id=ByKWUeWA->.
- 677 27. Li J, Cairns BJ, Li J, Zhu T. Generating synthetic mixed-type longitudinal electronic health  
678 records for artificial intelligent applications. *NPJ Digit Med.* 2023;6:98.
- 679 28. Xu L, Skoularidou M, Cuesta-Infante A, Veeramachaneni K. Modeling Tabular data using  
680 Conditional GAN. *arXiv [csLG]*. 2019.
- 681 29. Lederrey G, Hillel T, Bierlaire M. ciDATGAN: Conditional Inputs for Tabular GANs. *arXiv*  
682 *[csLG]*. 2022.
- 683 30. He Z, Tang X, Yang X, et al. Clinical Trial Generalizability Assessment in the Big Data Era:  
684 A Review. *Clin Transl Sci.* 2020;13:675–684.
- 685 31. Tripepi G, Jager KJ, Dekker FW, Zoccali C. Stratification for confounding--part 2: direct and  
686 indirect standardization. *Nephron Clin Pract.* 2010;116:c322-5.

- 687 32. Duan T, Rajpurkar P, Laird D, Ng AY, Basu S. Clinical Value of Predicting Individual  
688 Treatment Effects for Intensive Blood Pressure Therapy. *Circ Cardiovasc Qual Outcomes*.  
689 2019;12:e005010.
- 690 33. Brantner CL, Nguyen TQ, Tang T, Zhao C, Hong H, Stuart EA. Comparison of methods that  
691 combine multiple randomized trials to estimate heterogeneous treatment effects. *Stat Med*. 2024.  
692 Published online January 25, 2024. <https://doi.org/10.1002/sim.9955>.
- 693 34. Raghavan S, Josey K, Bahn G, et al. Generalizability of heterogeneous treatment effects  
694 based on causal forests applied to two randomized clinical trials of intensive glycemic control.  
695 *Ann Epidemiol*. 2022;65:101–108.
- 696 35. Fisher CK, Smith AM, Walsh JR, Coalition Against Major Diseases, Abbott, Alliance for  
697 Aging Research, Alzheimer’s Association, Alzheimer’s Foundation of America, AstraZeneca  
698 Pharmaceuticals LP, Bristol-Myers Squibb Company, Critical Path Institute, CHDI Foundation,  
699 Inc., Eli Lilly and Company, F. Hoffmann-La Roche Ltd, Forest Research Institute, Genentech,  
700 Inc., GlaxoSmithKline, Johnson & Johnson, National Health Council, Novartis Pharmaceuticals  
701 Corporation, Parkinson’s Action Network, Parkinson’s Disease Foundation, Pfizer, Inc., sanofi-  
702 aventis. Collaborating Organizations: Clinical Data Interchange Standards Consortium (CDISC),  
703 Epihian, Metrum Institute. Machine learning for comprehensive forecasting of Alzheimer’s  
704 Disease progression. *Sci Rep*. 2019;9:13622.
- 705 36. Walsh JR, Smith AM, Pouliot Y, Li-Bland D, Loukianov A, Fisher CK. Generating Digital  
706 Twins with Multiple Sclerosis Using Probabilistic Neural Networks. *arXiv [statML]*. 2020.

- 707 37. Bertolini D, Loukianov AD, Smith AM, et al. Modeling Disease Progression in Mild  
708 Cognitive Impairment and Alzheimer’s Disease with Digital Twins. *arXiv [csLG]*. 2020.
- 709 38. Eckardt J-N, Hahn W, Röllig C, et al. Mimicking clinical trials with synthetic acute myeloid  
710 leukemia patients using generative artificial intelligence. *NPJ Digit Med*. 2024;7:76.
- 711 39. Degtiar I, Rose S. A Review of Generalizability and Transportability. *Annual Review of*  
712 *Statistics and Its Application*. 2023;10:501–524.
- 713 40. Liu J, Li Y, Ge J, et al. Lowering systolic blood pressure to less than 120 mm Hg versus less  
714 than 140 mm Hg in patients with high cardiovascular risk with and without diabetes or previous  
715 stroke: an open-label, blinded-outcome, randomised trial. *Lancet*. 2024;404:245–255.
- 716 41. Center for Drug Evaluation and Research (CDER), Center for Biologics Evaluation and  
717 Research (CBER) U.S. Food and Drug Administration. Framework for FDA’s Real World  
718 Evidence Program. *US Food & Drug Administration*. 2018. Accessed March 6, 2024.  
719 <https://www.fda.gov/science-research/science-and-research-special-topics/real-world-evidence>.
- 720 42. Lederrey G, Hillel T, Bierlaire M. DATGAN: Integrating expert knowledge into deep  
721 learning for synthetic tabular data. *arXiv [csLG]*. 2022.
- 722 43. Khera R, Schuemie MJ, Lu Y, et al. Large-scale evidence generation and evaluation across a  
723 network of databases for type 2 diabetes mellitus (LEGEND-T2DM): a protocol for a series of  
724 multinational, real-world comparative cardiovascular effectiveness and safety studies. *BMJ*  
725 *Open*. 2022;12:e057977.



- 726 44. Goodfellow IJ, Pouget-Abadie J, Mirza M, et al. Generative Adversarial Networks. *arXiv*  
727 *[statML]*. 2014.
- 728 45. Hochreiter S, Schmidhuber J. Long short-term memory. *Neural Comput.* 1997;9:1735–1780.
- 729 46. Zhao Z, Kunar A, Van der Scheer H, Birke R, Chen LY. CTAB-GAN: Effective Table Data  
730 Synthesizing. *arXiv [csLG]*. 2021.
- 731 47. Arjovsky M, Chintala S, Bottou L. Wasserstein GAN. *arXiv [statML]*. 2017.
- 732 48. Normand ST, Landrum MB, Guadagnoli E, et al. Validating recommendations for coronary  
733 angiography following acute myocardial infarction in the elderly: a matched analysis using  
734 propensity scores. *J Clin Epidemiol.* 2001;54:387–398.
- 735 49. Patki N, Wedge R, Veeramachaneni K. The synthetic data vault. In: *2016 IEEE International*  
736 *Conference on Data Science and Advanced Analytics (DSAA)*. IEEE, 2016.
- 737 50. Zhao Z, Kunar A, Birke R, Chen LY. CTAB-GAN+: Enhancing Tabular Data Synthesis.  
738 *arXiv [csLG]*. 2022.
- 739 51. Kamthe S, Assefa S, Deisenroth M. Copula flows for synthetic data generation. *arXiv*  
740 *[statML]*. 2021.

741

742 **ACKNOWLEDGEMENTS**

743 The study is supported by the National Heart, Lung, and Blood Institute of the National Institutes  
744 of Health (R01HL167858). Dr. Thangaraj and Dr. Oikonomou are also supported by the National  
745 Heart, Lung, and Blood Institute of the National Institutes of Health (5T32HL155000-03 and  
746 1F32HL170592-01, respectively).

747

748 **CONTRIBUTIONS**

749 PMT and SVS contributed equally to the study. RK conceived the study and PMT, SVS, EKO,  
750 and RK drafted a research plan. EKO and PMT accessed and processed the data. PMT and SVS  
751 developed and analyzed the GAN model. PMT, SVS, EKO, and RK drafted the manuscript. All  
752 authors provided feedback regarding the study design and manuscript. RK supervised the study,  
753 procured funding, and is the guarantor.

754

755 **COMPETING INTERESTS**

756 The authors Dr. Thangaraj, Mr. Shankar, Dr. Oikonomou, and Dr. Khera are coinventors of a  
757 provisional patent related to the current work (63/606,203). Dr. Oikonomou is a co-inventor of  
758 the U.S. Patent Applications 63/508,315 63/177,117, a cofounder of Evidence2Health (with Dr.  
759 Khera), and has previously served as a consultant to Caristo Diagnostics Ltd (outside the present  
760 work). Dr. Nadkarni is a founder of Renalytix, Pensieve, and Verici and provides consultancy  
761 services to AstraZeneca, Reata, Renalytix, Siemens Healthineer, and Variant Bio, and serves a  
762 scientific advisory board member for Renalytix and Pensieve. He also has equity in Renalytix,  
763 Pensieve, and Verici. Dr. Mortazavi reported receiving grants from the National Institute of  
764 Biomedical Imaging and Bioengineering, National Heart, Lung, and Blood Institute, US Food

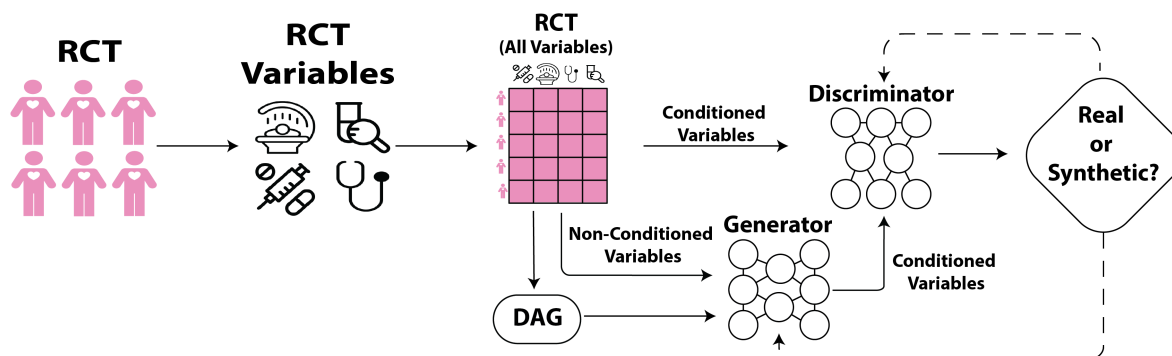
765 and Drug Administration, and the US Department of Defense Advanced Research Projects  
766 Agency outside the submitted work. In addition, B.J.M. has a pending patent on predictive  
767 models using electronic health records (US20180315507A1). Dr. Khera is an Associate Editor of  
768 JAMA. He receives support from the National Heart, Lung, and Blood Institute of the National  
769 Institutes of Health (under awards R01HL167858 and K23HL153775) and the Doris Duke  
770 Charitable Foundation (under award 2022060). He also receives research support, through Yale,  
771 from Bristol-Myers Squibb, Novo Nordisk, and BridgeBio. He is a coinventor of U.S. Pending  
772 Patent Applications 63/562,335, 63/177,117, 63/428,569, 63/346,610, 63/484,426, 63/508,315,  
773 and 63/606,203. He is a co-founder of Ensign-AI, Inc. and Evidence2Health, health platforms to  
774 improve cardiovascular diagnosis and evidence-based cardiovascular care.

775  
776  
777

778 FIGURES

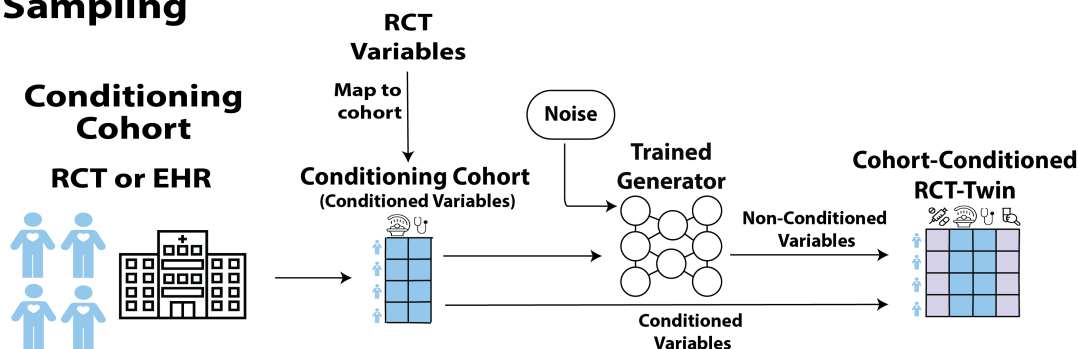
779 Figure 1: Graphical abstract of RCT-Twin-GAN model.

### A Training



780

### B Sampling

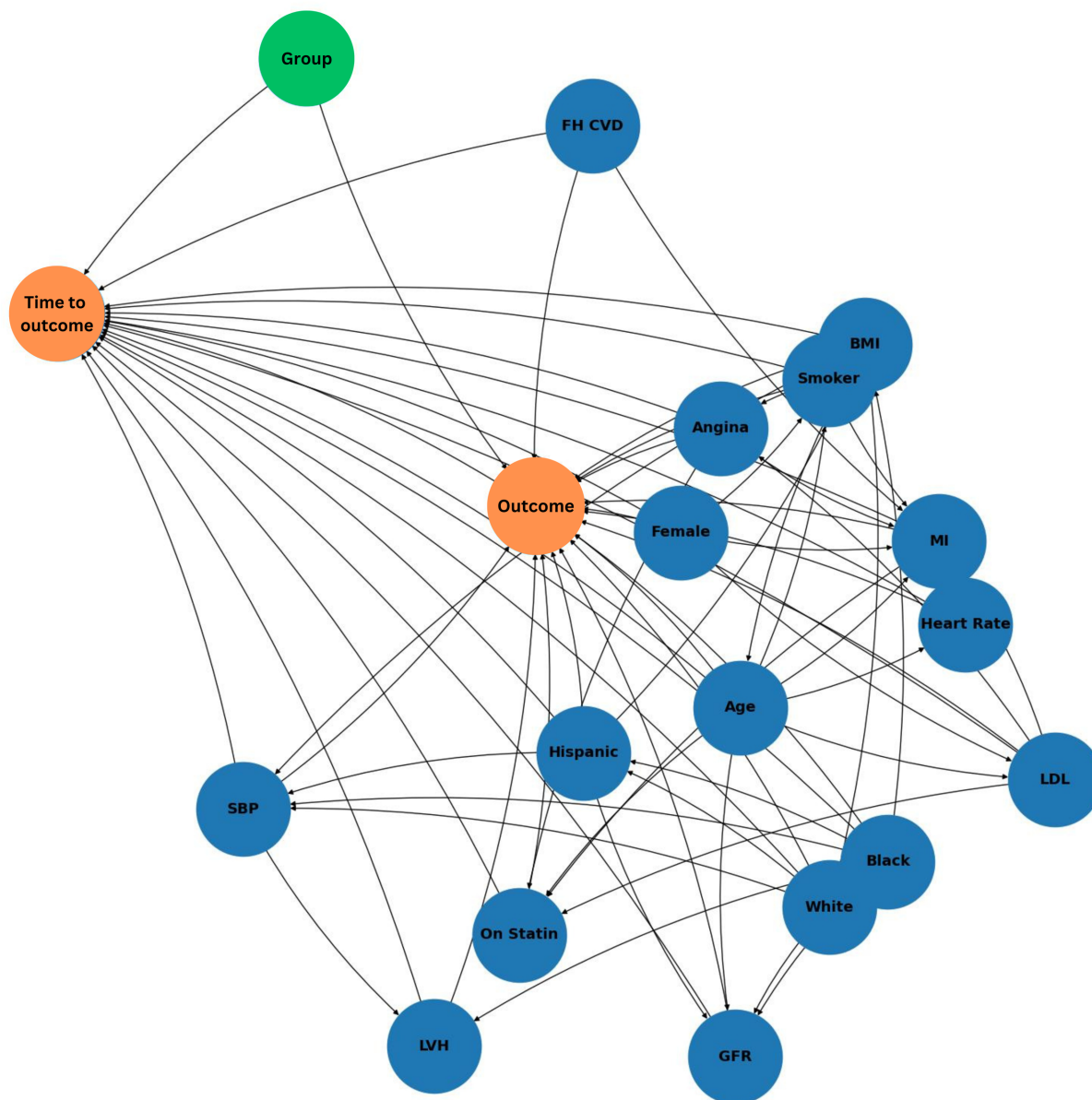


781

782 A. In the training phase, the original cohort (pink) is a randomized clinical trial (RCT), and  
 783 variables across all clinical domains are extracted from the cohort. The directed acyclic graph  
 784 (DAG) includes clinician-defined relationships between original cohort covariates and is inputted  
 785 to the generator, along with RCT values of the non-conditioned variables. The generator then  
 786 creates the conditioned variables, and the discriminator must differentiate from the original RCT  
 787 conditioned variables and the generator conditioned variables. Once the discriminator cannot  
 788 distinguish between the original and generated values, the training is complete. B. In the  
 789 sampling phase, conditioned variables from the RCT cohort are mapped to a conditioning cohort  
 790 (blue), examples of which are another RCT or a patient cohort in the electronic health record  
 791 (EHR). The trained generator then takes the conditioned variables from the conditioning cohort  
 792 and noise as input, and then generates non-conditioned variables. The final cohort-conditioned  
 793 RCT twin has conditioned covariate values from the conditioning cohort (blue) and generated  
 794 non-conditioned covariates based off the relationships and correlations between covariates (light  
 795 purple). Abbreviations: DAG: Directed Acyclic Graph, EHR: electronic health record, and RCT:  
 796 Randomized Clinical Trial.

797

798 **Figure 2: Directed acyclic graph of RCT-Twin-GAN.**  
799

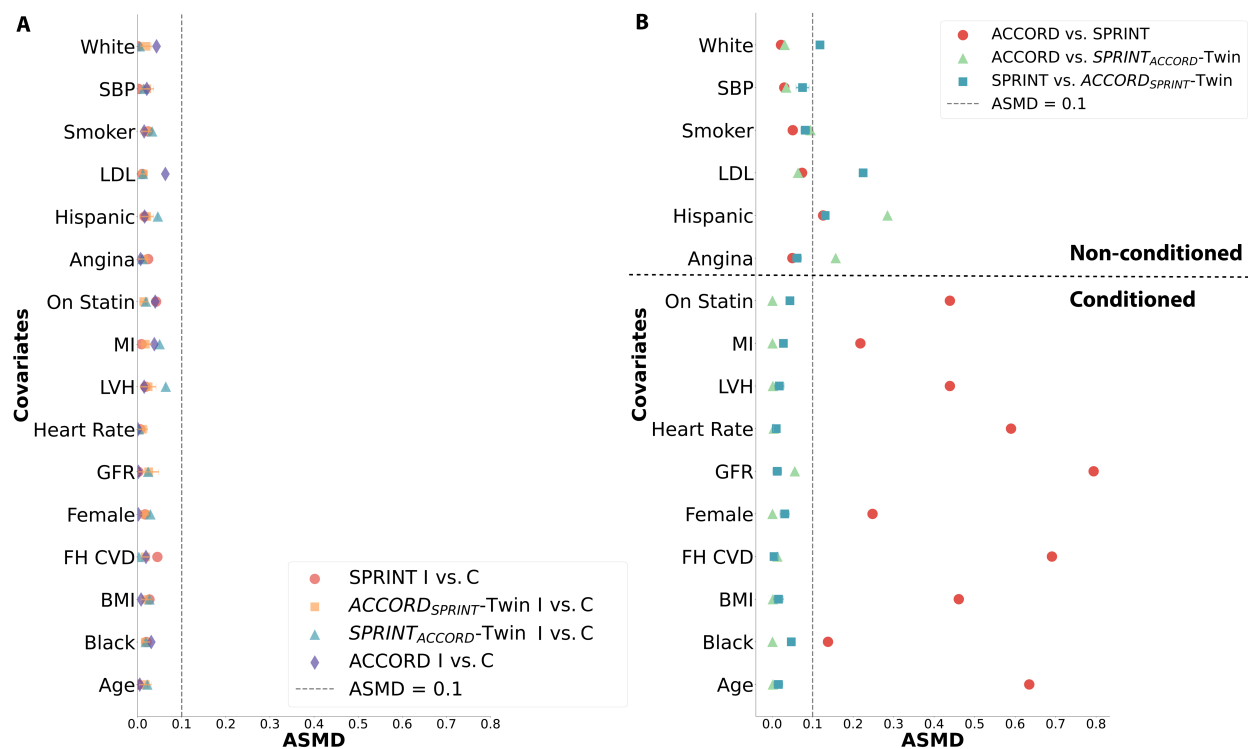


800 Directed relationships between covariates (in blue), time to outcome and outcome (in orange),  
801 and treatment arm designation (“Group”) in green. Abbreviations: BMI: Body Mass Index, CVD:  
802 Cardiovascular disease, FH: Family History, GFR: Glomerular Filtration Rate, LVH: Left  
803 ventricular hypertrophy, LDL: low-density lipoprotein, MI: Myocardial infarction, SBP: Systolic  
804 Blood Pressure.  
805  
806

807

808

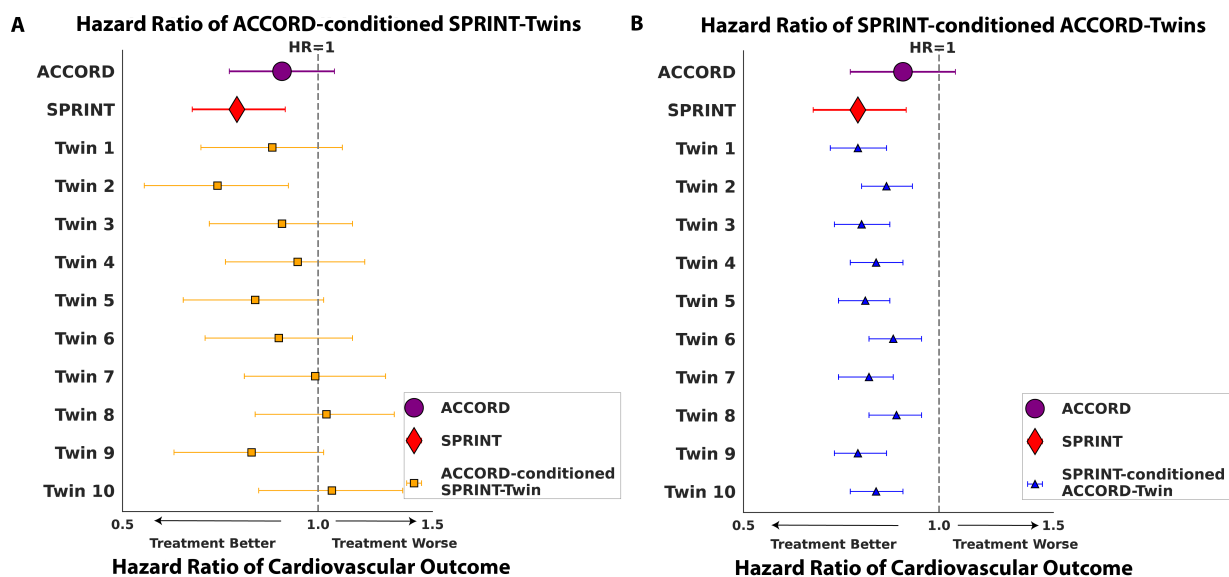
809 **Figure 3: Absolute standardized mean difference (ASMD) of covariates between datasets.**



810 (A) ASMD of covariates between treatment arms of RCTs and digital twins. Markers include  
 811 SPRINT (red circle), ACCORD<sub>SPRINT</sub> Twin (orange square), SPRINT<sub>ACCORD</sub>-Twin (blue  
 812 triangle), and ACCORD (purple diamond). The digital twin ASMDs are the mean of the 10  
 813 digital twin samples with standard deviation error bars. (B) ASMD of covariates between RCTs  
 814 and digital twins. Red circle represents ASMD between ACCORD and SPRINT, green triangle  
 815 represents ASMD between ACCORD and SPRINT<sub>ACCORD</sub>-Twins, and the blue square represents  
 816 the ASMD between SPRINT and ACCORD<sub>SPRINT</sub> Twins. The digital twin ASMDs are the mean  
 817 of the 10 digital twin samples with standard deviation error bars. The grey dotted line represents  
 818 an ASMD of 0.1, and the black dotted line separates non-conditioned and conditioned covariates.  
 819 The conditioning covariates included Age, Black, BMI, FH CVD, Female, GFR, Heart Rate,  
 820 LVH, MI, and On Statin. Abbreviations: ASMD: Absolute Standardized Mean Difference, BMI:  
 821 Body Mass Index, CVD: Cardiovascular disease, C: Control Arm, FH: Family History, eGFR:  
 822 Glomerular Filtration Rate, I: Intervention Arm, LDL: low-density lipoprotein, LVH: Left  
 823 ventricular hypertrophy, MI: Myocardial infarction, SBP: Systolic Blood Pressure.

825

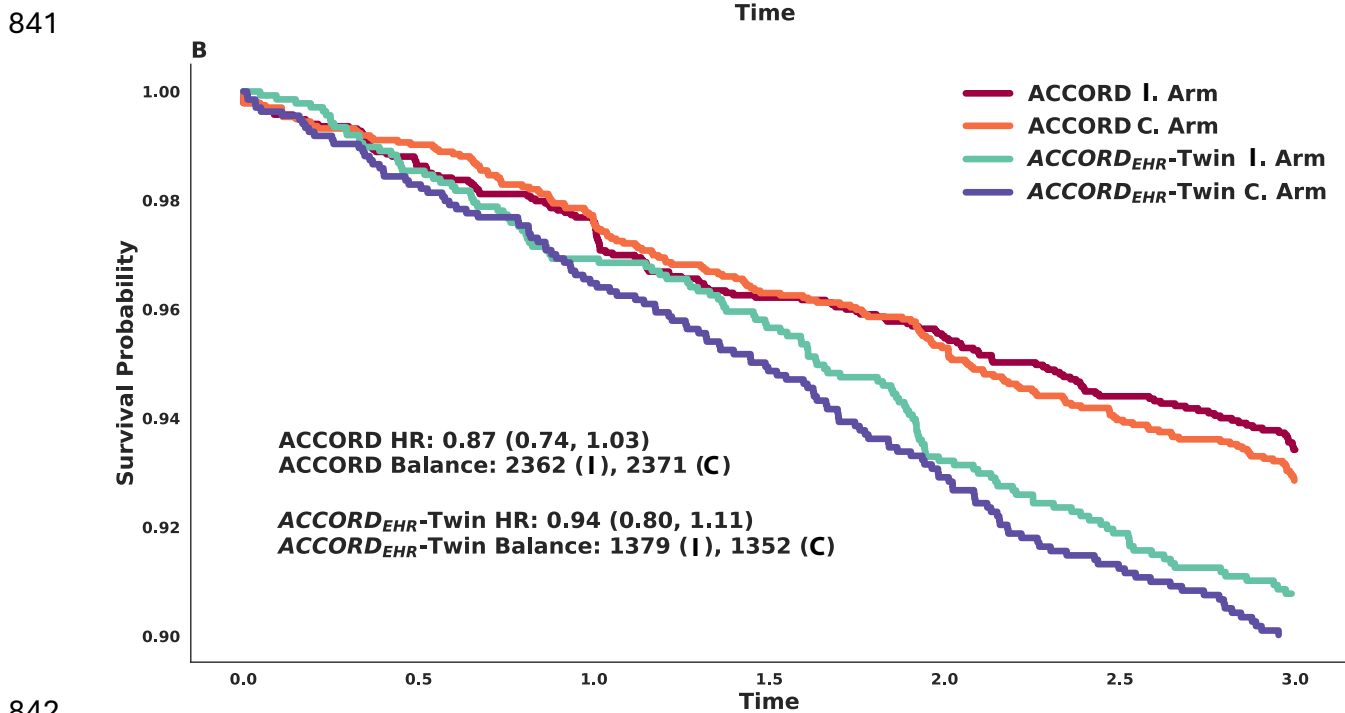
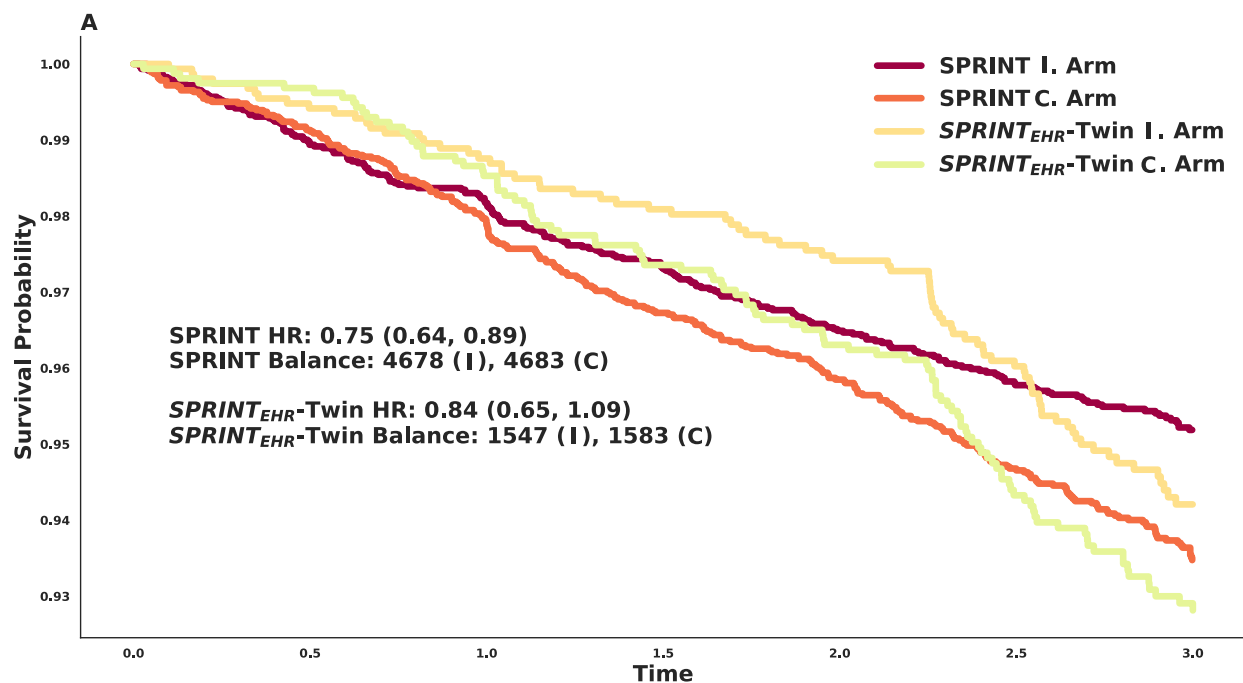
826 **Figure 4: Hazard ratios of intensive blood pressure lowering on cardiovascular outcomes.**  
827



828  
829 (A) Forest plot of the  $SPRINT_{ACCORD}$ -Twin datasets and (B) Forest plot of the  $ACCORD_{SPRINT}$ -  
830 Twin datasets. In both graphs, the purple circle is the ACCORD hazard ratio and 95% confidence  
831 interval, the red diamond is the SPRINT hazard ratio and 95% confidence interval, and the grey  
832 dotted line represents a hazard ratio of 1. In (A), orange squares are the hazard ratio of major  
833 cardiovascular outcome predicted for each twin run of ACCORD-conditioned SPRINT twins  
834 with 95% confidence intervals and in (B) the blue triangles are SPRINT-conditioned ACCORD  
835 twins. Abbreviations: MACE: Major Cardiovascular Outcomes,  $SPRINT_{ACCORD}$ -Twin:  
836 ACCORD-conditioned SPRINT Twin,  $ACCORD_{SPRINT}$ -Twin: SPRINT-conditioned ACCORD  
837 twin.

838

839 **Figure 5: Representative Kaplan-Meier curves of digital twins conditioned on EHR data.**  
840



842 (A) Kaplan-Meier curves of SPRINT treatment arms along with EHR-conditioned SPRINT  
843 treatment arms, (B) Kaplan Meier curves of ACCORD treatment arms along with EHR-  
844 conditioned ACCORD treatment arm balance of the original cohorts and digital twins.  
845 Abbreviations: ACCORD<sub>EHR</sub>Twin: ACCORD conditioned on EHR digital twin., C. Control arm,  
846 ACCORD: ACCORD, I: Intervention Arm, SPRINT<sub>EHR</sub>-Twin: SPRINT conditioned on  
847 EHR digital twin.  
848

849



850 **TABLES**

851 **Table 1: Minimal Requirements for RCT-Twin-GAN**

<b>Rules</b>
<ol style="list-style-type: none"><li>1. The model requires at least a pair of cohorts, one should be an RCT, the second should have the same covariates accessible and overlapping covariate distributions. This is, however, extendable to any number of cohorts.</li><li>2. There must be at least two treatment arms in the RCT with any ratio of randomization.</li><li>3. A measured outcome, categorical or continuous, should be available in the RCT to estimate treatment effect.</li><li>4. A sample size of at least 100 is needed for model convergence, but at least 1000 participants and further hyperparameter tuning of the model is needed to accurately estimate treatment effects.</li></ol>

852

853

854 **Table 2: Construction of the Directed Acyclic Graph**

1. Assemble covariates, treatment arms, time to outcome, and outcome.
2. Assign treatment arm as a source node that only connects to time to outcome and outcome.
3. Ensure all covariates are connected to time to outcome and outcome.
4. Clinicians assess which covariates influence each other based off clinical knowledge, such as systolic blood pressure to left ventricular hypertrophy or current smoker to angina or MI.
5. Pearson and Spearman correlations are calculated between every unconnected pair of variables and those with both Pearson and Spearman correlations  $>0.75$  are assessed.
6. Clinicians assess which of the suggested pairs are clinically relevant and add them to the DAG.
7. Repeat steps 5 and 6 until no more clinically relevant pairs are suggested.

855

856

857 Table 3: CiDATGAN Architecture, Training and Sampling  
858

**1. DAG Construction:** For each cohort (e.g., SPRINT and ACCORD), we constructed a DAG representing the causal relationships between covariates and outcomes. This DAG was used to inform the generator about the relevant correlation structure and prevent overfitting of correlations from noise.

The graph  $G$  for a DAG is specified by the modeler to define the correlations between the variables in the data.

- Each variable  $V_t$  in the table  $T$  must be associated with a node in the graph  $G$ .
- A directed edge between two nodes, i.e.  $V_{t1} \rightarrow V_{t2}$ , means that the generation of the first variable  $V_{t1}$  will influence the generation of the second variable  $V_{t2}$ .
- The absence of a link between two variables means that their correlation is not directly learned by the Generator.

Once the DAG  $G$  was created, we defined several sets:

- $A(V_t)$ : the set of ancestors of the variable  $V_t$
- $D(V_t)$ : the set of direct ancestors of the variable  $V_t$
- $S(V_t)$ : the set of source nodes leading to the variable  $V_t$
- $E(V_t)$ : the set of in-edges of the variable  $V_t$

**2. Training and Conditioning Phase:** The DATGAN architecture was trained on the original cohort data (e.g., SPRINT). Continuous and categorical variables from the dataset were encoded at the discriminator input. We introduced conditioning covariates to condition the generation of synthetic data. The DAG was constructed so that the conditioning covariates were removed and treated as source nodes because the discriminator was only trained on non-conditioned covariates. The generator utilized Gaussian noise, the DAG-informed covariate relationships, and the conditioned covariate values from the conditioning cohort to generate synthetic data. The discriminator was trained to distinguish between the generated synthetic data and the real cohort data.

The mathematical representations of these are included below:

$T_0$ : Original cohort with  $N_0$  variables ( $v_i^0$  for  $i = 1, \dots, N_0$ )

$T_C$ : Original cohort with  $N_C$  variables ( $v_i^C$  for  $i = 1, \dots, N_C$ ) where  $N_0 > N_C$

$N_{CV}$ : Number of common variables across both cohorts such that  $N_{CV} \leq N_0$

These common variables are denoted as  $T_0^{ci}$  and  $T_C^{ci}$  for the original and conditioning cohorts respectively.

Goal: To generate complementary variables  $T_0^c = T_0 - T_0^{ci}$  using the values of common variables  $T_C^{ci}$  as inputs.

**Generator:**

Let,

$G$  = Generator

$z$  = Gaussian Noise

$T_C^{ci}$  = Conditional Inputs

Generates each variable in  $T_0^c$

Let  $T_{O \rightarrow C}^{c, synth}$  = Generated Data

### Discriminator:

Let,

D = Discriminator

D distinguishes between the real data  $T_0^c$  and generated data  $T_{O \rightarrow C}^{c, synth}$

The model is trained on  $T_0$  with the associated DAG structure where the conditional variables are source nodes.

The generation of synthetic variables  $V_t^O$  in  $T_0^c$  using LSTM cells follows an order provided by the linearization of the DAG.

### LSTM Cells:

Each LSTM cell  $LSTM_t$  is associated with the variable  $v_t^O$ , ordered based on the DAG. The cell takes as input the cell state of the previous variable in the DAG  $C_{t-1}$  and the input tensor  $i_t$ , which is a concatenation of:

$$i_t = [z_t, f_{t-1}, a_t]$$

Where:

- $z_t$  is a tensor of Gaussian noise.
- $f_{t-1}$  is the transformed output of the previous LSTM cell in the DAG.
- $a_t$  is the attention vector used to retain information from previous ancestors not directly linked to the current cell in the DAG.

### Attention Vector $a_t$ :

$$a_t = \sum_{k \in A(t) \setminus P(t)} \frac{\exp(\alpha_k^t)}{\sum_{j=1}^{|A(t)|} \exp(\alpha_j^t)} f_k$$

where  $A(t) \setminus P(t)$  is the set of ancestors of the variable  $v_t^C$  in the DAG, excluding direct predecessors,  $\alpha^t$  is a learned attention weight vector, and  $f_k$  is the final output of the LSTM cell  $LSTM_k$ .

### Output of LSTM Cells:

Each LSTM cell outputs two tensors, the new cell state  $C_t$  and the output of the cell  $h_t$ . This output is then passed through fully connected layers to get the synthetic values  $v_t^{C, synth}$ :

$$v_t^{O, synth} = FC(h_t)$$

The synthetic tensor is resized to a common size between all variables using an input transformer.

### Handling Conditional Inputs:

For variables in  $T_0^{ci}$ , the generator needs the transformed output  $f_{t-1}$  of the direct ancestor and the direct output  $h_k$  of all ancestors.

Transformed Output  $f_t$ :

The same type of input transformer is used to get  $f_t$  for the conditional inputs.

**Dense (Fully connected) Layer Transformation:**

Since LSTM output  $h_t$  is not available for conditional inputs, the original value  $v_t$  is transformed using a Dense layer:

$$h_t = \text{Dense}(v_t)$$

The parameters in this Dense layer are learned during the training process, allowing the model to use the conditional inputs in the attention vector.

**3. Sampling Process:**

During the sampling phase, the model received Gaussian noise, the DAG, and the values of conditioned covariates for the patient, which can come from either the original dataset  $T_O^{ci}$  or the conditioning dataset  $T_C^{ci}$ . This combination was used to generate the final digital twin ( $T_{O \rightarrow C}^{c, synth}$ ) from each patient of the conditioning cohort, which included a copy of conditioned covariate values, and generated non-conditioned covariates based off the correlations between the covariates from the original cohort.

**4. Sampling and Iterations:**

Digital twin generation was repeated for a specified number of iterations (e.g., 10 iterations). In each iteration, the generator produced a synthetic dataset based on the conditioned covariates and Gaussian noise. Each iteration generated a complete digital twin cohort.”

859