

1 Assessing racial bias in type 2 diabetes risk prediction algorithms

2
3 Hélène T. Cronjé, PhD^{1,¶}, Alexandros Katsiferis, MSc^{1,¶}, Leonie K. Elsenburg, PhD¹, Thea O.
4 Andersen, MSc¹, Naja H. Rod, PhD¹, Tri-Long Nguyen, PhD¹, Tibor V. Varga, PhD^{1,*}

5
6 Short title: *Racial bias in diabetes risk scores*

7
8 Affiliations:

- 9 1. Section of Epidemiology, Department of Public Health, University of Copenhagen,
10 Copenhagen, Denmark

11
12 [¶] These authors contributed equally

13 * Corresponding author:

14 Tibor V. Varga

15 ORCID: <https://orcid.org/0000-0002-2383-699X>

16 Postal address:

17 Bartholinsgade 6Q; DK-1356 Copenhagen K; Denmark

18 Tel: +45 35 32 77 39

19 E-mail: tibor.varga@sund.ku.dk

20

21 Abstract

22 Risk prediction models for type 2 diabetes can be useful for the early detection of individuals at
23 high risk. However, models may also bias clinical decision-making processes, for instance by
24 differential risk miscalibration across racial groups. We investigated whether the Prediabetes
25 Risk Test (PRT) issued by the National Diabetes Prevention Program, and two prognostic
26 models, the Framingham Offspring Risk Score, and the ARIC Model, demonstrate racial bias
27 between non-Hispanic Whites and non-Hispanic Blacks. We used National Health and Nutrition
28 Examination Survey (NHANES) data, sampled in six independent two-year batches between
29 1999 and 2010. A total of 9,987 adults without a prior diagnosis of diabetes and with fasting
30 blood samples available were included. We calculated race- and year-specific average predicted
31 risks of type 2 diabetes according to the risk models. We compared the predicted risks with
32 observed ones extracted from the US Diabetes Surveillance System across racial groups
33 (summary calibration). All investigated models were found to be miscalibrated with regards to
34 race, consistently across the survey years. The Framingham Offspring Risk Score overestimated
35 type 2 diabetes risk for non-Hispanic Whites and underestimated risk for non-Hispanic Blacks.
36 The PRT and the ARIC models overestimated risk for both races, but more so for non-Hispanic
37 Whites. The risk of type 2 diabetes extracted from these landmark models were more severely
38 overestimated for non-Hispanic Whites compared to non-Hispanic Blacks, potentially resulting
39 in a larger fraction of non-Hispanic Whites being prioritized for a preventive intervention, but
40 also more likely to be overdiagnosed and overtreated, with a larger fraction of non-Hispanic
41 Blacks being potentially underprioritized and undertreated.

42

43 Introduction

44 Early detection of individuals who are at high risk for developing type 2 diabetes is a powerful
45 strategy to tackle the diabetes epidemic through targeted prevention[1]. Type 2 diabetes affects
46 one in eight adults in the United States (US) and is the nation's seventh leading cause of death. A
47 look beyond the averages reveals a disparity in the incidence and mortality rate of diabetes
48 across racial groups, with non-Hispanic Whites being the only group wherein rates are below the
49 national average[2, 3].

50 Despite their comparatively lower risk[4], non-Hispanic White groups remain overrepresented in
51 diabetes risk prediction literature[5]. Consequently, the implementation of evidence-based risk
52 prediction in primary care is vulnerable to limited generalizability to other racial groups. Biased
53 prediction models may prioritize individuals of certain racial groups for preventive action at
54 different rates, or at different stages in their disease progression[6].

55 Today, the current standards of care focus on identifying individuals with prevalent type 2
56 diabetes or prediabetes in asymptomatic adults and referring them to appropriate prevention or
57 treatment routes [10]. Screening relies mainly on the Prediabetes Risk Test (PRT) [10,11], a
58 diagnostic model developed by the National Diabetes Prevention Program
59 (<https://diabetes.org/diabetes/risk-test> and <https://www.cdc.gov/prediabetes/risktest>). To the best
60 of our knowledge, despite the availability of numerous other US-based *prognostic* risk prediction
61 models [7-9], very few or none of them are currently implemented in routine.

62 In this project, we set out to investigate whether the PRT[7] and two key US-based prognostic
63 prediction models for type 2 diabetes that informed the current guidelines, demonstrate
64 consistent (mis)calibration across Black and White non-Hispanic US residents. We further

65 investigate to what extent the current clinical guidelines may address or perpetuate these biases
66 and discuss concrete action points for stakeholders to (further) promote fair prognostic disease
67 prediction.

68 Research Design and Methods

69 Study population

70 The Continuous National Health and Nutrition Examination Survey (NHANES) is a repeated
71 cross-sectional survey that has been providing a vast amount of publicly available health data on
72 a randomly selected nationally representative sample of the civilian, non-institutionalized US
73 population since 1999[8]. Demographic, questionnaire, physical examination, and biochemical
74 data are collected in two-year intervals. Reliability of sub-group data is ensured by the
75 oversampling of minority racial groups and persons over the age of 60.

76 Our population of interest comprised non-Hispanic White and non-Hispanic Black adults at risk
77 of type 2 diabetes. We included NHANES survey content between 1999 and 2010, sampled in
78 six independent 2-year batches (1999-2000, 2001-2002, ..., 2009-2010). A total of 62,160
79 participants were surveyed during this timeframe, of whom 22,661 provided fasting blood
80 samples. Analyses were restricted to the latter group due to the utilization of fasting biochemical
81 measures in the prediction algorithms tested. Upon excluding individuals with diagnosed type 2
82 diabetes (based on self-report and medication use), those younger than 18 years, and those who
83 did not self-report their race as non-Hispanic White or non-Hispanic Black (n=12,674), we
84 retained an analytical sample of 9,987 individuals (**Figure 1, Appendix Figure 1**). We used
85 fasting subsample weights to correct for representativeness in our analyses and allow for
86 inferences made from this subset to be representative of the larger US population of interest.

87

88 Prediction models

89 Our analyses focused on the PRT[7] issued by the National Diabetes Prevention Program, and
90 two widely known North-American prognostic type 2 diabetes risk models: the Framingham
91 Offspring Risk Score[9] and the Atherosclerosis Risk in Communities (ARIC) Risk Model[10]
92 (**Appendix Table 1**). We used the NHANES database to extract the data needed to calculate
93 model estimates according to the authors' guidelines (**Figure 1**). Continuous predictors included
94 age (years), body mass index (BMI, kg/m²), waist circumference (cm), height (cm), systolic and
95 diastolic blood pressure (mmHg), fasting glucose (mmol/L), high-density lipoprotein cholesterol
96 (HDL-C) (mmol/L), and triglycerides (mmol/L).

97 Binary predictors were sex (male/female), race (non-Hispanic Black and non-Hispanic White),
98 anti-hypertensive medication use (yes/no), previously diagnosed hypertension (yes/no), physical
99 activity (yes/no, coded yes for datasets between 1999-2006 if the participant responded “More
100 active” to a single question “Compared with most men/women your age, would you say that you
101 are...”, and coded yes for datasets between 1999-2006 if the participant reported either
102 occupational or leisure time vigorous physical activity), and family history of diabetes (yes/no;
103 yes when participants answered in the affirmative to any of the specific family members
104 questioned explicitly between the 1999 and 2004 surveys, and from 2005 when participants
105 answered yes to the single question of whether a close relative has diabetes). For the
106 Framingham Offspring Risk Score and ARIC Risk Model, we included the latter-described
107 variable as a proxy for *parental history of diabetes* due to the unavailability of a better-defined
108 variable in NHANES.

109

110 Statistical analysis

111 All analyses were undertaken using R 4.1.2. The NHANES data were extracted and prepared for
112 analysis using the framework depicted in **Appendix Figure 1**. Missing data were imputed using
113 multivariate imputation by chained equations[11]. For all variables, random forest was utilized
114 (five iterations), and 15 imputed copies were generated. The imputation models included the
115 sampling weights to account for the complex survey design[12]. Convergence was visually
116 inspected for randomly selected imputed datasets. Race- and survey-specific estimates derived
117 from the final analytical cohort (N=9,987) from the multiple imputation framework were pooled
118 using Rubin's rules[13].

119 We used the Framingham Offspring Risk Score and the ARIC Risk Model to predict the risk of
120 type 2 diabetes of participants. Given those predictions, we subsequently calculated race- and
121 year-specific average predicted incidence proportions and 95% confidence intervals (CIs) for
122 eight or nine years (depending on the model used) subsequent to each survey batch interval.
123 Those race- and year- specific averages were weighted to account for the survey design of the
124 NHANES data. We also determined the proportion of each survey batch that would be identified
125 for further screening based on their PRT score. The latter assigns a score between -1 and 9, with
126 values of 5 or higher indicating a high risk of either prediabetes or type 2 diabetes. Cumulative
127 type 2 diabetes incidences (matched to the respective risk score time frames) by race were
128 calculated as the sum of yearly reported incidences extracted from the US Diabetes Surveillance
129 System database[14]. We subsequently compared the latter with the NHANES-derived weighted
130 average predicted risks of the two prognostic models. For illustration purposes we also plotted
131 the observed cumulative incidences against the predicted proportion of undiagnosed type 2
132 diabetes cases of the PRT test (**Figure 1**). Thus, we evaluated the predictive performance of the

133 model in terms of overall calibration within different racial subgroups, by calculating expected-
134 to-observed incidence risk ratios[15].

135

136 Ethics statement

137 The NHANES was subject to ethics review by the National Center for Health Statistics Research

138 Ethics Review Board and approval was obtained bi-annually between each full proposal review.

139 All participants provided informed consent.

140

141 Role of the funding source

142 This project was supported by the University of Copenhagen, the Novo Nordisk Foundation, and

143 the Independent Research Fund Denmark. The funders had no role in the study design; in the

144 collection, analysis, and interpretation of data; in the writing of the report; and in the decision to

145 submit the paper for publication.

146 Results

147 Race-stratified weighted, unimputed descriptive statistics of the study population are presented
148 per survey year in **Table 1**. Imputed descriptive statistics are shown in **Appendix Table 2**.

149 Compared to non-Hispanic White NHANES participants, non-Hispanic Black groups were
150 younger, less likely to have resided in the US from birth or to have obtained a high school
151 diploma, and more likely to be hypertensive or have a family history of diabetes. While non-
152 Hispanic Blacks were also more likely to have higher BMIs than their White counterparts, they
153 had more favorable triglyceride profiles, and did not differ from non-Hispanic White groups in
154 terms of waist circumference or HDL-C. For survey years spanning 1999-2004, non-Hispanic
155 Whites were proportionally less physically active than non-Hispanic Blacks, with the trend
156 reversing from 2005 onward.

157 **Appendix Table 3** presents the race-stratified age-adjusted type 2 diabetes incidence rates for
158 non-Hispanic Whites and non-Hispanic Blacks from the US Diabetes Surveillance System[14].
159 Overall, type 2 diabetes incidence peaked in 2008 and has been decreasing since, although
160 incidences remain higher than what they were in 2000. There continues to be considerable
161 variation across racial groups, with non-Hispanic Blacks consistently having higher incidence
162 rates compared to non-Hispanic Whites.

163 Predicted average risks for each examined type 2 diabetes model were compared with
164 cumulative incidences calculated from the US Diabetes Surveillance System. As an example, we
165 used the NHANES cohort from 1999 to calculate race-stratified 8-year average predicted type 2
166 diabetes risk, and thus, calculated that by 2007, 7% of non-Hispanic Whites, and 6% of non-
167 Hispanic Blacks are expected to develop type 2 diabetes. We compared our calculated estimates
168 with real-life race-stratified cumulative incidences from the matching 8-year period (2000-2007)

169 (cumulative incidences for this period are 5% for non-Hispanic Whites and 8% for non-Hispanic
170 Blacks). We repeated this analysis for the three models, and each of the six cohorts from
171 NHANES between 1999 and 2010, to obtain race-stratified predicted estimates until 2017
172 (**Appendix Table 4**). Results for the three models are shown in **Figure 2**.

173 While the Framingham Offspring Risk Score overestimated type 2 diabetes risk for non-Hispanic
174 Whites, it underestimated risk for non-Hispanic Blacks. The PRT and the ARIC Model
175 overestimated type 2 diabetes risk for both races, but more so for non-Hispanic Whites compared
176 to non-Hispanic Blacks. The PRT showed the highest overestimation; based on the scores
177 approximately half of the total populations of the cohorts were prioritized for screening.

178 We calculated the ratios of the calculated average predicted risks (incidence proportion) to the
179 true cumulative incidences (**Figure 3**). All three models demonstrated overestimation of risk for
180 non-Hispanic Whites. The ARIC model delivered a lower overestimation for non-Hispanic
181 Blacks (average ratio = 1.22) compared to non-Hispanic Whites (average ratio = 2.31). Similarly,
182 the PRT delivered a lower overestimation for non-Hispanic Blacks (average ratio = 5.05)
183 compared to non-Hispanic Whites (average ratio = 9.51) when considering 8-year risk of
184 diabetes.

185 We visualized the correspondence between the dichotomized PRT scores of individuals and the
186 individual predicted probabilities from the two prognostic type 2 diabetes risk prediction
187 algorithms (**Figure 4**). The PRT scores showed moderate correlation with the predicted
188 probabilities by the Framingham (Kendall's tau=0.39) and the ARIC models (Kendall's
189 tau=0.56). Above and below the threshold of scoring five on the PRT, non-Hispanic Whites
190 demonstrated a slightly higher risk of type 2 diabetes compared to non-Hispanic Blacks using

191 both the Framingham and the ARIC models, however this difference did not reach statistical
192 significance (95% CIs overlapping 0).

193

194

195 Discussion

196 In recent decades, algorithmic decision making has become a routine part of healthcare. Simple
197 risk scores are routinely used to evaluate an individual's risk of developing a disease for most
198 cardiometabolic outcomes, including type 2 diabetes[16], and complex artificial intelligence-
199 driven models are continuously developed to predict common complications of type 2
200 diabetes[17]. Nonetheless, regardless of how powerful artificial intelligence models can be in
201 capturing complex interactions and patterns in data, the appropriateness of the available datasets
202 in terms of representativeness and quality, remains crucial to algorithmic design.

203 Recent examples have shown that models developed within biased datasets, or not directly
204 addressing data biases, will likely propagate inequalities into clinical decisions[6, 18, 19]. As a
205 result of biased algorithmic decision making, those who are already marginalized and less
206 represented, may encounter further obstacles in accessing optimal healthcare, generating a
207 vicious cycle. In this report, we investigated whether the PRT—a nationally adopted screening
208 algorithm for prediabetes and type 2 diabetes—and two landmark type 2 diabetes predictive
209 models developed in the US were racially biased. We compared the proportion of the population
210 identified as being at high risk for incident type 2 diabetes with corresponding reported national
211 statistics on type 2 diabetes incidence by racial group, within the predicted time horizon.

212 We interpret the model-derived proportions, indicating the population at risk for type 2 diabetes
213 within a certain timeframe, as the fraction of the racial groups that could be prioritized for
214 preventive intervention in that timeframe. When models underestimate the proportion of the
215 population at risk for type 2 diabetes compared with national statistics, that indicates that the
216 models identify less individuals in need of preventive action than the actual needs based on
217 national statistics. Conversely, in case of overestimation, larger fractions of the population could

218 be prioritized for preventive action and larger proportions would potentially receive support to
219 achieve their health targets than would be necessary based on national statistics. Imbalances in
220 over- and underestimation between racial groups reflect algorithmic bias, i.e., systematic
221 differences in the performance of algorithms across specific groups. Overdiagnosis, and
222 overtreatment pose significant challenges, such as the detection and unnecessary treatment of
223 individuals who will ultimately remain asymptomatic[20]. However, in the case of type 2
224 diabetes, the consensus on preventive action largely involves non-invasive lifestyle interventions,
225 which are less likely to pose harm to false positives compared to e.g., pharmaceutical treatment.
226 Nonetheless, imposing individuals to unnecessary, even nonaggressive, procedures, can lead to
227 negative psychosocial consequences, i.e., increased stress levels, reduced quality of life, or social
228 stigma. Ultimately, and given the nature and guidelines of the specific chronic disease, we
229 consider a potential underestimation of risk as a bigger threat compared to overestimation.

230 Our results show that PRT and the two examined prognostic type 2 diabetes risk prediction
231 models consistently demonstrated a larger risk overestimation for non-Hispanic Whites than
232 Blacks. Moreover, the Framingham Offspring Risk Score, developed in a 99% White and non-
233 Hispanic population, underestimated type 2 diabetes risk for non-Hispanic Blacks.

234 Although the nationally adopted PRT appears to be a highly sensitive (i.e., very unlikely to
235 misclassify diabetes cases as non-cases), it is also non-specific, and identifies approximately half
236 of the total population for the screening of type 2 diabetes, very likely falsely identifying
237 individuals who are normoglycemic. Despite the wide net the PRT casts, this risk scoring still
238 identifies proportionally more non-Hispanic Whites for preventive action than non-Hispanic
239 Blacks. When compared to the observed national incidence rates, the PRT fails to fully account

240 for the racial differences that we observe in real-life incidence rates. However, a certain amount
241 of overestimation from PRT was expected, given that it also indicates cases of prediabetes.
242 Notably, all three examined algorithms are almost exclusively accounting for inherent risk
243 (biological sex and family history of diabetes) or markers of metabolic health (biochemical
244 profiles, blood pressure, and anthropometry). The ARIC Risk Model was the only algorithm to
245 include race, a non-metabolic parameter; and the PRT was the only algorithm incorporating a
246 lifestyle factor (physical activity). In this nationally representative dataset, the three examined
247 models consistently predicted higher average risks for non-Hispanic Whites compared to non-
248 Hispanic Blacks, contrary to the official national statistics. This observation indicates that the
249 metabolic health variables implemented in the models failed to capture the observed lower risk
250 of type 2 diabetes in non-Hispanic Whites compared to Blacks. On the contrary, we observed
251 higher predicted risks for non-Hispanic Whites than for non-Hispanic Blacks. While our results
252 are only suggestive, such phenomenon could be explained by the absence in the models of socio-
253 economic related determinants describing health literacy status, access to healthcare, the latter
254 being on average higher among non-Hispanic Whites, thus serving as compensatory, protective
255 factors that would result in lower realized risks of type 2 diabetes compared to non-Hispanic
256 Blacks. Our descriptive statistics of the analyzed data confirm these patterns, with non-Hispanic
257 Blacks having lower educational attainment on average and a higher likelihood of being
258 immigrants. As an acute solution to the observed algorithmic bias resulting from these health
259 inequalities, the examined algorithms – and specifically the PRT that is already adopted by
260 healthcare – would likely benefit from the explicit inclusion of (additional) markers related to
261 education, health-literacy, and other socioeconomic determinants (that are expected to correlate
262 with race)[21, 22]. The inclusion of such features in the models would likely shift the

263 distributions of predicted probabilities higher for non-Hispanic Blacks and thus prioritize a larger
264 fraction of this population for preventive action.

265 Given the underrepresentation of minorities in studies where models were developed, limited
266 algorithmic generalizability might be expected within some racial groups. When it comes to
267 sample size considerations in developing novel algorithms, three major approaches are to be
268 considered. The first approach is to develop models in nationally representative populations. This
269 will result in predictive models that will perform more accurately for the majority. The cohort
270 investigated by the developers of the ARIC Model[10] most closely represented this approach by
271 comprising 85% non-Hispanic Whites and 15% non-Hispanic Blacks. In the original publication,
272 the ARIC model did perform most accurately in the majority population, although the positive
273 coefficient in the risk equation for Blacks increased the risk estimates for this group resulting in
274 an accuracy nearing that of the majority population. When tested externally using NHANES
275 data, however, we still see marked differences in model performance between non-Hispanic
276 Whites and Blacks, reflecting algorithmic bias. The second approach is to develop models in
277 cohorts with roughly equal sample sizes across groups, resulting in models that are expected to
278 perform with the same confidence and precision across groups, but may perform less optimally
279 for the majority. The third approach is to develop separate models in separate groups. This
280 approach was taken when developing the Reynolds Risk Scores separately for men[23] and
281 women[24] to predict sex-specific 10-year cardiovascular disease risk. While both total sample
282 size and sample size of respective groups are important considerations in the development of
283 prediction models, societal and data biases (e.g. differential access to healthcare) may also
284 impact the development and subsequent performance of algorithms, as observed before[6].

285 Prediction models can perpetuate and reinforce data biases that lie in the core of the societies we
286 live in. Explicit counteraction is warranted to address these biases when developing prediction
287 models. In recent years, algorithmic fairness, as an emergent research area, has been receiving
288 increasing attention and the development of toolkits for the development of fair models have
289 taken flight[25]. It has been raised that without the consideration of sensitive attributes, it is
290 cumbersome to correct for data biases, and develop fair models[26]. However, the inclusion of
291 race as a feature in predictive models has been at the center of academic debate[27-29]. For
292 instance, the vast majority of clinical prediction models for cardiovascular diseases do not
293 include race in the assessment of individual risk[27]. Arguments for the omission of these
294 sensitive attributes are that racial discrimination can potentially be reinforced with the addition
295 of race in the models (i.e. racial profiling), and that there is only weak evidence for genetic or
296 biological differences across races[30]. Additionally, researchers have expressed concern that
297 race insertion might create risks of falsely interpreting racial inequities as immutable facts rather
298 than disparities that require intervention[31, 32]. Nonetheless, variables encoding race are
299 deemed to reflect a complex combination of factors related to biological, cultural, as well as
300 socio-contextual aspects. Therefore, while the use of race as a predictor may suffer from crucial
301 flaws, it can also pave the way for better algorithmic decision-making across various racial
302 groups, for example by promoting personalized care, targeting those in the most need of
303 intervention[27]. The exclusion of race from algorithms may lead to differential external validity
304 and potentially harmful decision making, as models may fail to account for differential risk
305 across groups[30]. Undeniably, we are facing a ‘data problem’: data currently used to generate
306 models are likely to embed inherent racial imbalances, lack intersectionality, and/or important
307 social-contextual features that potentially confound the association between race and various

308 clinical endpoints[29, 31]. As an example, Obermeyer *et al.* demonstrated significant racial bias
309 in a predictive algorithm using healthcare expenditure as predictive features for prioritizing
310 individuals for healthcare interventions[6]. Thus, it is crucial for future studies to investigate
311 complex diseases through a ‘structural racism lens’[33, 34] and disentangle causal relationships
312 between various socioeconomic factors, human physiology, and outcomes, in pursuit of
313 identifying the most appropriate candidate features for predictive modeling.

314 Based on these findings, we suggest that any published and/or candidate diagnostic or prognostic
315 models should demonstrate algorithmic fairness before adoption in healthcare, e.g., via a
316 systematic comparison of their performance in external samples stratified by sensitive attributes,
317 such as race, ethnicity, and sex. Until fairness is a default consideration in algorithmic
318 development, and models adopted by healthcare are thoroughly checked for fairness, it is likely
319 that societal inequalities will keep propagating into clinical decisions. Thus, we call for medical
320 journals to recommend reporting criteria related to the fairness of novel predictive algorithms.

321 While numerous toolkits, checklists, and governance frameworks have been developed on the
322 topic[35], we emphasize five key points that journals may want to require reporting on: (i)
323 potential data biases; (ii) possible measurements of bias; (iii) utilized bias mitigation strategies;
324 (iv) potential impacts of the biases; (v) and the generalizability of results across relevant
325 sensitive attributes. We recommend that expert groups and healthcare policymakers responsible
326 for the adoption of any risk models (e.g., cutoffs, simple risk scores, or complex algorithms) in
327 clinical practice consider fairness a decisive factor. In the specific context of type 2 diabetes,
328 there is a need to develop an accurate risk score for new-onset diabetes that helps identify groups
329 with unmet needs within minorities.

330 Our study needs to be considered under some limitations. First, the predictive performance of the
331 models could not be assessed in terms of discrimination as in their traditional reported form. The
332 repeated cross-sectional design of the NHANES did not allow us to compare the individual
333 predicted probabilities returned by the prediction models with the individual observed outcomes,
334 since participants were not followed-up over time. While prediction models are often used for
335 decision-making at the individual level, one needs to bear in mind that predicted probabilities are
336 in fact drawn from *groups* of persons with comparable characteristics in terms of predictor
337 values. To show disparities of predictive performance across racial groups, we reported ratios of
338 average predicted to average observed type 2 diabetes incidences as an overall, summary
339 measure of calibration within the different racial groups[15]. The interpretation of our findings
340 thus remains at the group level. Second, we caution against an interpretation of the trend over
341 time in the estimated predictive performance, as changes in observed incidence over time could
342 be merely due to changes in case detection. Skewed case detection across racial groups (i.e.,
343 differential measurement error), could exaggerate or attenuate the reported disparities. This is,
344 however, less of a concern due to the lack of observable differential detection in the available
345 national diagnosed/undiagnosed type 2 diabetes data[36]. Third, the complex survey design of
346 NHANES complicated the estimation procedure in the presence of missing item responses. To
347 address this issue, we performed multiple imputation, while including the sampling weights in
348 the imputation models[12], and using random forest models to allow for complex interactions
349 and non-linearities to be captured[37]. Last, we acknowledge that a large number of type 2
350 diabetes models are available, and it would be possible to extend our list. Here we selected
351 widely known models developed in the US, but future reports could possibly include a wider
352 range of models for testing.

353 In summary, our study shows that the PRT currently adopted by US healthcare, and prognostic
354 type 2 diabetes prediction models available for adoption in US healthcare are likely attached with
355 some degree of racial bias, which in turn is likely to perpetuate inequalities by providing fewer
356 benefits to minorities, who already demonstrate higher risk for metabolic diseases. We provide
357 specific recommendations that we believe can improve the standards of publishing and the
358 adoption of algorithmic processes in healthcare.

359

360 Additional Information

361 Author contributions

362 Authors' contributions: The corresponding authors attest that all listed authors meet the ICMJE
363 authorship criteria and that no others meeting the criteria have been omitted. AK, TOA, TN, and
364 TVV performed the statistical analysis. HTC, AK, LKE, TOA, TN, and TVV drafted the
365 manuscript. HTC, AK, and TVV created the visualizations. HTC, LKE performed a literature
366 review for the project. TN, and TVV were responsible for the conceptualization of the project.
367 TVV was responsible for project management. TVV is the guarantor of this manuscript. All
368 listed authors interpreted the results, reviewed, and edited the manuscript, verified the data, and
369 approved the final, submitted version.

370

371 Conflict-of-interest statement

372 The authors declare no conflict of interest.

373

374 Data availability statement

375 NHANES data is publicly available and can be accessed via
376 <https://www.cdc.gov/nchs/nhanes/index.htm>

377 The authors are willing to share codes.

378

379

380 Funding

381 HTC and AK are supported by a grant from Novo Nordisk Foundation Challenge Programme for
382 the project Harnessing the Power of Big Data to Address the Societal Challenge of Aging
383 (NNF17OC0027812). TOA is supported by a grant from the Independent Research Fund
384 Denmark (7025-00005B). TVV is supported by the "Data Science Investigator - Emerging 2022"
385 grant from Novo Nordisk Foundation (NNF22OC0075284) and the Department of Public Health
386 (University of Copenhagen).

387 References

- 388 1. Dunkley AJ, Bodicoat DH, Greaves CJ, Russell C, Yates T, Davies MJ, et al. Diabetes
389 prevention in the real world: effectiveness of pragmatic lifestyle interventions for the prevention
390 of type 2 diabetes and of the impact of adherence to guideline recommendations: a systematic
391 review and meta-analysis. *Diabetes care*. 2014;37(4):922-33.
- 392 2. Control CfD, Prevention. National diabetes statistics report, 2020. Atlanta, GA: Centers
393 for Disease Control and Prevention, US Department of Health and Human Services. 2020:12-5.
- 394 3. Heron M. Deaths: Leading Causes for 2019. *Natl Vital Stat Rep*. 2021;70(9):1-114. Epub
395 2021/09/15. PubMed PMID: 34520342.
- 396 4. Wang MC, Shah NS, Carnethon MR, O'Brien MJ, Khan SS. Age at diagnosis of diabetes
397 by race and ethnicity in the United States from 2011 to 2018. *JAMA internal medicine*.
398 2021;181(11):1537-9.
- 399 5. Ayensa-Vazquez JA, Leiva A, Tauler P, López-González AA, Aguiló A, Tomás-Salvá
400 M, et al. Agreement between Type 2 Diabetes Risk Scales in a Caucasian Population: A
401 Systematic Review and Report. *Journal of clinical medicine*. 2020;9(5):1546.
- 402 6. Obermeyer Z, Powers B, Vogeli C, Mullainathan S. Dissecting racial bias in an algorithm
403 used to manage the health of populations. *Science*. 2019;366(6464):447-53.
- 404 7. Bang H, Edwards AM, Bomback AS, Ballantyne CM, Brillon D, Callahan MA, et al.
405 Development and validation of a patient self-assessment score for diabetes risk. *Annals of*
406 *internal medicine*. 2009;151(11):775-83.
- 407 8. Statistics. CfDCaPNCfH. National Health and Nutrition Examination Survey. [
408 <https://www.cdc.gov/nchs/nhanes/index.htm>]. Accessed: 26 Jan 2022.
- 409 9. Wilson PW, Meigs JB, Sullivan L, Fox CS, Nathan DM, D'Agostino RB. Prediction of
410 incident diabetes mellitus in middle-aged adults: the Framingham Offspring Study. *Archives of*
411 *internal medicine*. 2007;167(10):1068-74.
- 412 10. Schmidt MI, Duncan BB, Bang H, Pankow JS, Ballantyne CM, Golden SH, et al.
413 Identifying individuals at high risk for diabetes: The Atherosclerosis Risk in Communities study.
414 *Diabetes care*. 2005;28(8):2013-8.
- 415 11. Van Buuren S, Groothuis-Oudshoorn K. mice: Multivariate imputation by chained
416 equations in R. *Journal of statistical software*. 2011;45:1-67.
- 417 12. Kim JK, Michael Brick J, Fuller WA, Kalton G. On the bias of the multiple - imputation
418 variance estimator in survey sampling. *Journal of the Royal Statistical Society: Series B*
419 *(Statistical Methodology)*. 2006;68(3):509-21.
- 420 13. Rubin D. *Multiple Imputation for Nonresponse in Surveys* Wiley J&Sons New York.
421 1987.
- 422 14. US Diabetes Surveillance System. [
423 <https://gis.cdc.gov/grasp/diabetes/DiabetesAtlas.html#>]. Accessed: 26 Jan 2022.
- 424 15. Debray TP, Damen JA, Snell KI, Ensor J, Hooft L, Reitsma JB, et al. A guide to
425 systematic review and meta-analysis of prediction model performance. *Bmj*. 2017;356.
- 426 16. Association AD. <https://www.diabetes.org/risk-test>.
- 427 17. Gulshan V, Peng L, Coram M, Stumpe MC, Wu D, Narayanaswamy A, et al.
428 Development and Validation of a Deep Learning Algorithm for Detection of Diabetic
429 Retinopathy in Retinal Fundus Photographs. *JAMA : the journal of the American Medical*

- 430 Association. 2016;316(22):2402-10. Epub 2016/11/30. doi: 10.1001/jama.2016.17216. PubMed
431 PMID: 27898976.
- 432 18. Sjoding MW, Dickson RP, Iwashyna TJ, Gay SE, Valley TS. Racial bias in pulse
433 oximetry measurement. *New England Journal of Medicine*. 2020;383(25):2477-8.
- 434 19. Chapman CH, Schechter CB, Cadham CJ, Trentham-Dietz A, Gangnon RE, Jagsi R, et
435 al. Identifying equitable screening mammography strategies for Black women in the United
436 States using simulation modeling. *Annals of internal medicine*. 2021;174(12):1637-46.
- 437 20. Vogt H, Green S, Ekstrøm CT, Brodersen J. How precision medicine and screening with
438 big data could increase overdiagnosis. *Bmj*. 2019;366.
- 439 21. Laraia BA, Karter AJ, Warton EM, Schillinger D, Moffet HH, Adler N. Place matters:
440 neighborhood deprivation and cardiometabolic risk factors in the Diabetes Study of Northern
441 California (DISTANCE). *Social science & medicine*. 2012;74(7):1082-90.
- 442 22. Kimenai DM, Pirondini L, Gregson J, Prieto D, Pocock SJ, Perel P, et al. Socioeconomic
443 Deprivation: An Important, Largely Unrecognized Risk Factor in Primary Prevention of
444 Cardiovascular Disease. *Circulation*. 2022;146(3):240-8.
- 445 23. Ridker PM, Paynter NP, Rifai N, Gaziano JM, Cook NR. C-reactive protein and parental
446 history improve global cardiovascular risk prediction: the Reynolds Risk Score for men.
447 *Circulation*. 2008;118(22):2243-51.
- 448 24. Ridker PM, Buring JE, Rifai N, Cook NR. Development and validation of improved
449 algorithms for the assessment of global cardiovascular risk in women: the Reynolds Risk Score.
450 *JAMA : the journal of the American Medical Association*. 2007;297(6):611-9.
- 451 25. Bellamy RK, Dey K, Hind M, Hoffman SC, Houde S, Kannan K, et al. AI Fairness 360:
452 An extensible toolkit for detecting and mitigating algorithmic bias. *IBM Journal of Research and
453 Development*. 2019;63(4/5):4: 1-4: 15.
- 454 26. Zliobaite I. Fairness-aware machine learning: a perspective. arXiv preprint
455 arXiv:170800754. 2017.
- 456 27. Paulus JK, Kent DM. Race and ethnicity: a part of the equation for personalized clinical
457 decision making? *Circulation: Cardiovascular Quality and Outcomes*. 2017;10(7):e003823.
- 458 28. Essien UR, Jackson LR. Race effects in CVD prediction models. *Journal of general
459 internal medicine*. 2019;34(4):484-.
- 460 29. Waters EA, Colditz GA, Davis KL. Essentialism and Exclusion: Racism in Cancer Risk
461 Prediction Models. *JNCI: Journal of the National Cancer Institute*. 2021.
- 462 30. Paulus JK, Wessler BS, Lundquist CM, Kent DM. Effects of race are rarely included in
463 clinical prediction models for cardiovascular disease. *Journal of general internal medicine*.
464 2018;33(9):1429-30.
- 465 31. Vyas DA, Eisenstein LG, Jones DS. Hidden in plain sight—reconsidering the use of race
466 correction in clinical algorithms. *Mass Medical Soc*; 2020. p. 874-82.
- 467 32. Lett E, Asabor E, Beltrán S, Cannon AM, Arah OA. Conceptualizing, Contextualizing,
468 and Operationalizing Race in Quantitative Health Sciences Research. *Annals of family
469 medicine*.2792.
- 470 33. Adkins-Jackson PB, Chantarat T, Bailey ZD, Ponce NA. Measuring structural racism: a
471 guide for epidemiologists and other health researchers. *American journal of epidemiology*. 2021.
- 472 34. Robinson WR, Renson A, Naimi AI. Teaching yourself about structural racism will
473 improve your machine learning. *Biostatistics*. 2020;21(2):339-44.

- 474 35. Madaio MA, Stark L, Wortman Vaughan J, Wallach H, editors. Co-designing checklists
475 to understand organizational challenges and opportunities around fairness in ai. Proceedings of
476 the 2020 CHI Conference on Human Factors in Computing Systems; 2020.
- 477 36. Prevention CfDca. Prevalence of Both Diagnosed and Undiagnosed Diabetes. [
478 <https://www.cdc.gov/diabetes/data/statistics-report/diagnosed-undiagnosed-diabetes.html>].
479 Accessed: 26 Jan 2022.
- 480 37. Shah AD, Bartlett JW, Carpenter J, Nicholas O, Hemingway H. Comparison of random
481 forest and parametric imputation models for imputing missing data using MICE: a CALIBER
482 study. American journal of epidemiology. 2014;179(6):764-74.
- 483

484 Tables

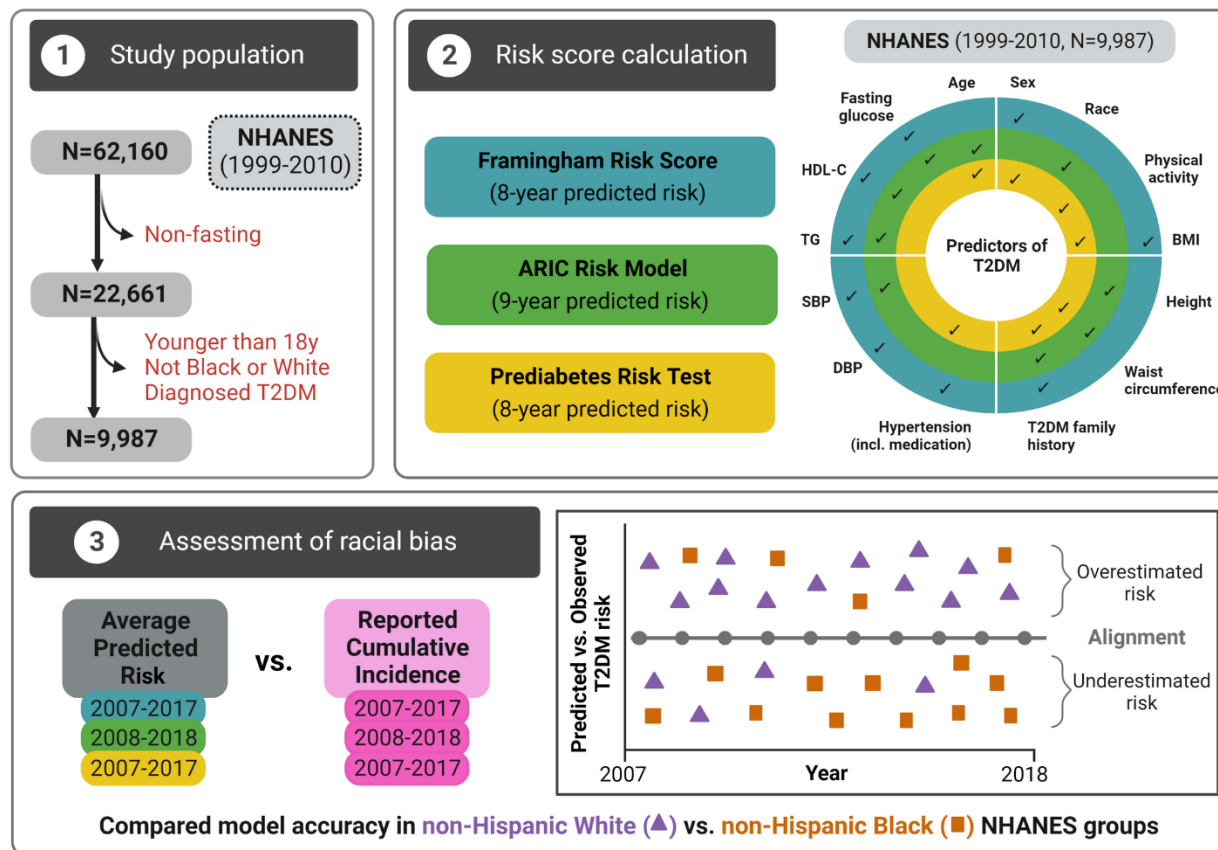
485 **Table 1.** Descriptive statistics of the unimputed NHANES data (N=9,987).

	1999–2000		2001–2002		2003–2004		2005–2006		2007–2008		2009–2010	
	non-Hispanic White	non-Hispanic Black	non-Hispanic White	non-Hispanic Black	non-Hispanic White	non-Hispanic Black	non-Hispanic White	non-Hispanic Black	non-Hispanic White	non-Hispanic Black	non-Hispanic White	non-Hispanic Black
N	1,222	175	1,575	231	1,451	221	1,450	221	1,474	224	1,502	241
Age, years	44.9 (17.1)	40.2 (14.5)	45.4 (16.9)	39.9 (15.5)	45.6 (17.2)	40.1 (16.8)	46.4 (17.5)	40.7 (16.2)	46.5 (17.3)	40.8 (15.6)	47.1 (17.5)	41.1 (15.5)
Female, %	52.2	54.2	52.0	55.0	52.1	55.6	50.9	54.8	51.7	54.4	52.0	54.3
Highschool graduate, %	82.8	62.1	84	69.2	84.6	66.1	86.5	71.7	83.9	70.6	85.8	71.9
Missing, %	3.5	4.9	3.5	5.1	4.5	6.1	3.4	7.5	3.5	6.1	3.0	4.5
Born in the US, %	95.0	86.6	95.8	92.4	94.4	93.9	95.5	90.8	96.1	92.0	94.0	84.5
Physically active, %	37.5	40.7	35.7	37.9	34.0	38.9	38.2	37.1	43.2	42.5	40.1	37.2
Family history of diabetes, %	44.5	46.4	46.5	46.9	45.5	52.4	34.7	41.9	31.4	37.9	29.8	42.7
Missing, %	5.3	5.2	5.3	6.7	6.7	7.5	6.5	8.6	5.6	6.9	4.8	5.3
Body mass index, kg/m ²	26.9 (5.8)	29.1 (7.1)	27.7 (6.1)	28.8 (7.6)	27.8 (6.1)	29.9 (7.6)	28.2 (7.0)	30.3 (7.8)	27.6 (5.8)	29.1 (7.1)	28.2 (6.4)	30.5 (8.1)
Missing, %	0.5	1.8	6.6	3.6	1.4	2.0	1.7	2.2	1.9	2.1	1.2	0.4
Waist circumference, cm	93.5 (15.4)	94.9 (16.6)	95.5 (15.4)	93.9 (16.7)	96.9 (15.6)	97.1 (17.0)	97.0 (16.6)	98.0 (17.5)	96.3 (15.1)	96.2 (16.2)	97.8 (16.0)	98.1 (16.9)
Missing, %	1.7	4.1	5.2	4.2	4.6	6.7	4.5	4.9	6.6	5.4	4.4	5.2
Height, cm	170.3 (9.9)	169.9 (9.5)	170.1 (9.9)	169.6 (9.9)	170.8 (9.9)	169.7 (9.4)	170.4 (9.9)	169.6 (9.7)	170.8 (10.1)	170.5 (9.0)	170.5 (9.7)	169.0 (9.7)
Missing, %	0.5	1.8	3.6	1.2	1.3	2.0	1.6	2.0	1.8	2.1	0.9	0.4
Fasting glucose, mmol/L	5.1 (1.0)	5.2 (1.5)	5.1 (0.9)	5.1 (1.0)	5.2 (0.7)	5.2 (0.9)	5.2 (1.0)	5.1 (0.8)	5.1 (0.7)	5.1 (1.1)	5.2 (0.8)	5.2 (0.8)
Missing, %	5.6	10.4	4.4	13.5	5.5	6.1	5.3	9.4	4.1	14.7	5.1	9.7
HDL-C, mmol/L	1.3 (0.4)	1.4 (0.4)	1.3 (0.4)	1.4 (0.4)	1.4 (0.4)	1.5 (0.4)	1.5 (0.4)	1.5 (0.5)	1.4 (0.4)	1.6 (0.4)	1.4 (0.4)	1.4 (0.5)
Missing, %	5.7	10.8	4.5	14.1	5.2	5.9	5.1	8.8	4.2	14.1	4.9	9.5
Triglycerides, mmol/L	1.6 (1.1)	1.1 (0.7)	1.7 (2.1)	1.2 (0.7)	1.7 (1.4)	1.2 (0.8)	1.6 (1.2)	1.2 (1.2)	1.5 (1.1)	1.0 (0.6)	1.4 (1.1)	1.1 (0.7)
Triglycerides, ln(mmol/L)	0.3 (0.5)	0.0 (0.5)	0.3 (0.6)	0.0 (0.5)	0.3 (0.6)	0.0 (0.5)	0.3 (0.6)	0.0 (0.5)	0.2 (0.5)	-0.1 (0.6)	0.2 (0.5)	0.0 (0.5)
Missing, %	5.9	11.0	4.5	14.1	5.2	6.3	5.3	9.9	4.2	14.1	4.9	9.5
Systolic BP, mmHg	121.2 (17.7)	123.2 (17.8)	122.4 (18.8)	124.1 (19.4)	121.2 (17.9)	124.8 (20.7)	121.2 (17.2)	125.3 (18.0)	119.2 (16.0)	121.5 (17.3)	118.2 (15.5)	122.9 (16.7)
Diastolic BP, mmHg	72.0 (11.9)	74.1 (12.2)	72.6 (11.5)	73.4 (13.2)	70.4 (13.1)	71.2 (12.8)	68.8 (12.4)	69.8 (14.7)	69.3 (11.6)	70.4 (12.5)	67.8 (11.7)	71.1 (13.5)
Missing, %	2.4	6.5	3.5	5.9	5.1	5.5	3.8	8.5	4.6	6.2	3.6	7.6
Antihypertensive use, %	14.9	15.4	16.1	19.1	18.4	18.4	18.5	21.8	21.2	22.4	21.3	21.4
Missing, %	27.1	28.6	28.2	28.6	24.8	24.8	22.2	21.5	18.2	17.3	17.5	22.1
Hypertension diagnosis, %	41.4	44.3	43.7	46.9	41.4	41.4	39.7	44.2	37.1	41.5	37.0	44.4

486

487 Figure legends

488 **Figure 1.** Study design

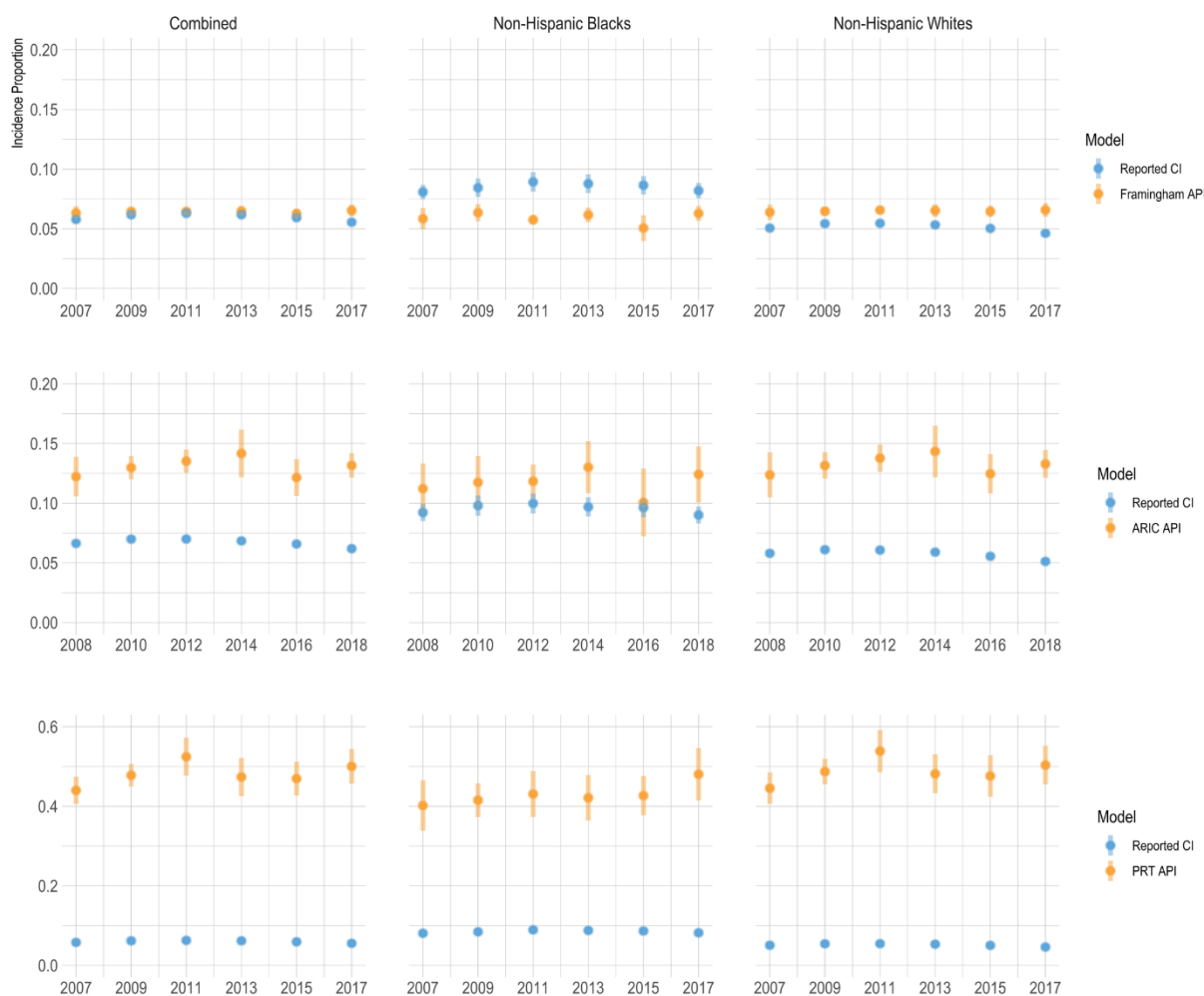


489

490

491

492 **Figure 2.** The first and third panel show the 8-year predicted type 2 diabetes risk for the
493 Framingham Offspring Risk Score and the Prediabetes Risk Test and 8-year cumulative incidences
494 of type 2 diabetes from the US Diabetes Surveillance System, combined and per racial group. The
495 second panel shows the 9-year predicted type 2 diabetes risks for the ARIC Model and 9-year
496 cumulative incidences of type 2 diabetes from the US Diabetes Surveillance System, combined
497 and per racial group.



498
499 Abbreviations: API - average predicted incidence (proportion); CI - cumulative incidence; PRT –
500 Prediabetes Risk Test.
501 Error bars reflect 95% confidence intervals.

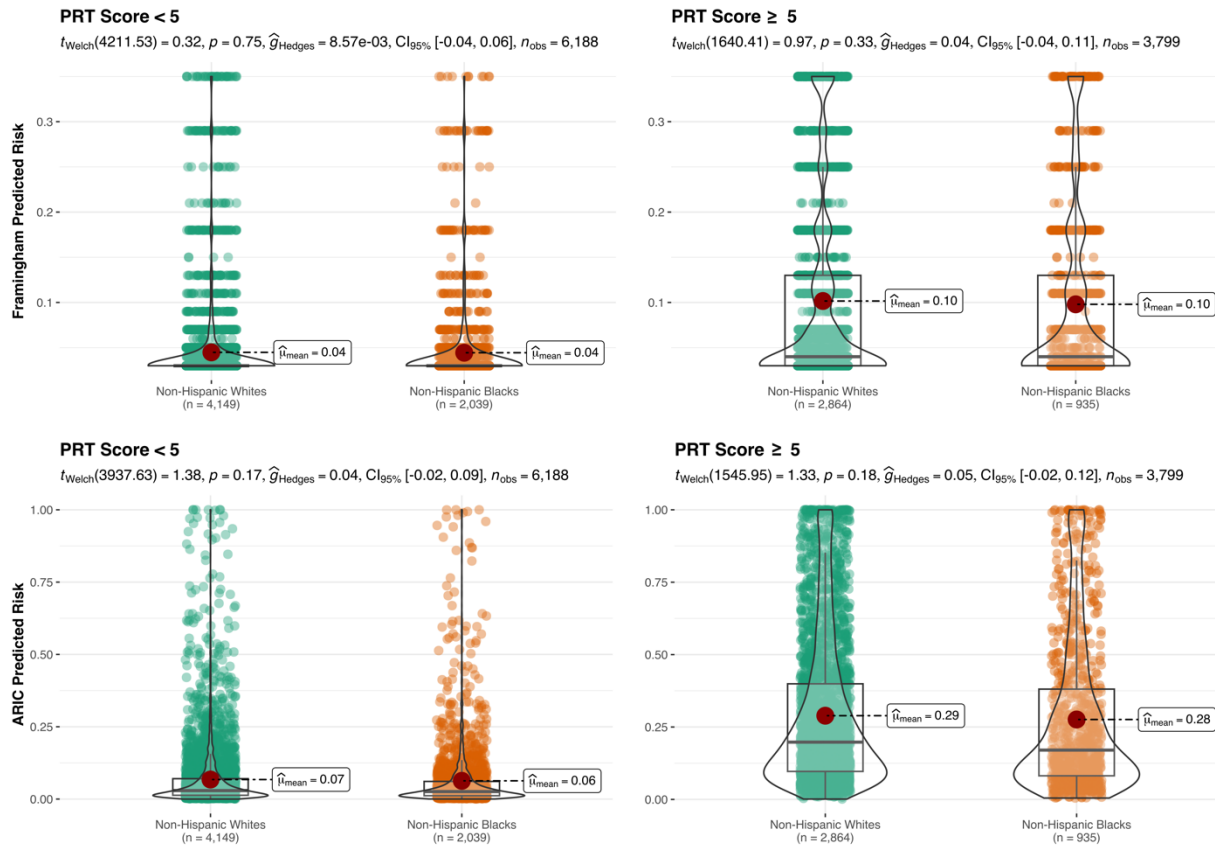
502 **Figure 3.** Ratios between predicted Framingham Offspring Risk Score, ARIC model, and
503 Prediabetes Risk Test incidence proportions and reported cumulative incidences of type 2 diabetes
504 overall, and per racial group.



505
506 Abbreviations: API - average predicted incidence (proportion); CI - cumulative incidence.
507

508 **Figure 4.** Box-Violin Plots of Framingham Offspring Risk Score and ARIC Model individual
509 predicted risks in groups scoring less or greater than 5 in the Prediabetes Risk Test.

510



511

512 Group differences were estimated using Welch's t-tests.

513

514 Appendix headings

515 Appendix Figure 1: Study flowchart

516 Appendix Table 1. Type 2 diabetes risk prediction model included in the analytic framework.

517 Appendix Table 2. Descriptive statistics of the imputed NHANES data (N=9,987).

518 Appendix Table 3. Age-adjusted incidence rates of type 2 diabetes per 1000 individuals.

519 Appendix Table 4. Predicted type 2 diabetes risks per racial group by the three prediction models, and

520 cumulative incidences of type 2 diabetes from the US Diabetes Surveillance System.

521

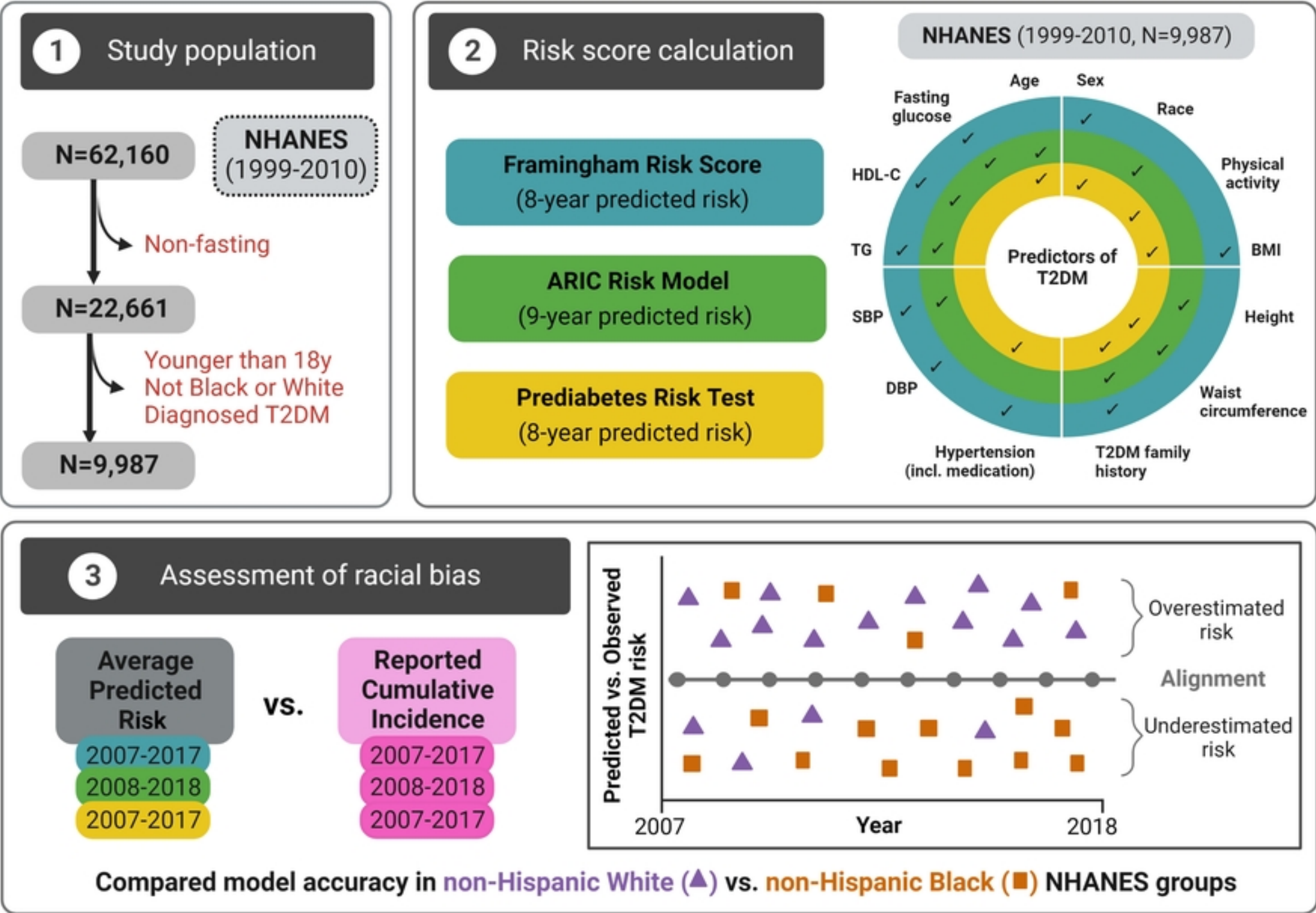


Figure 1

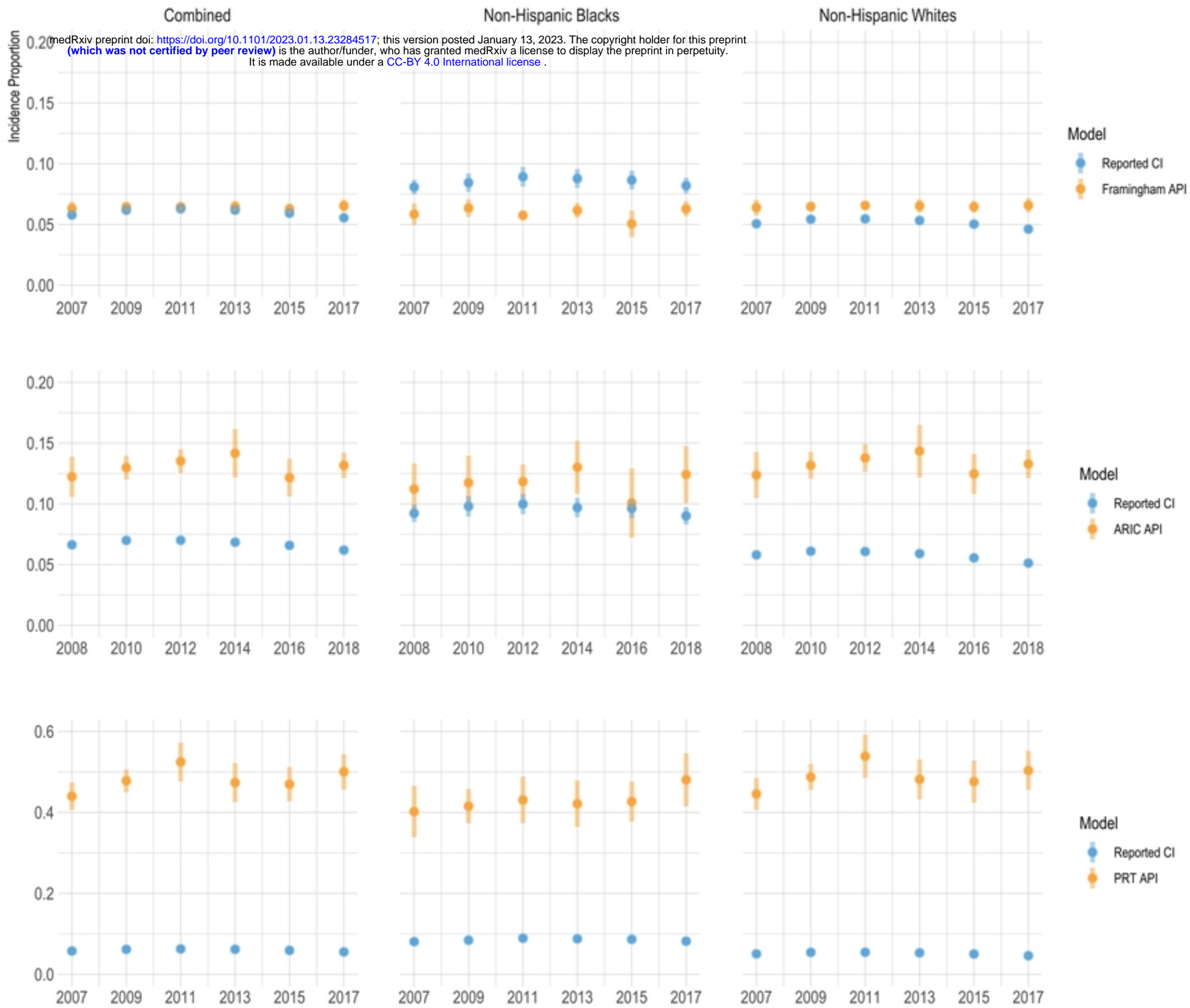


Figure 2

Race ● Combined ● Non-Hispanic Blacks ● Non-Hispanic Whites

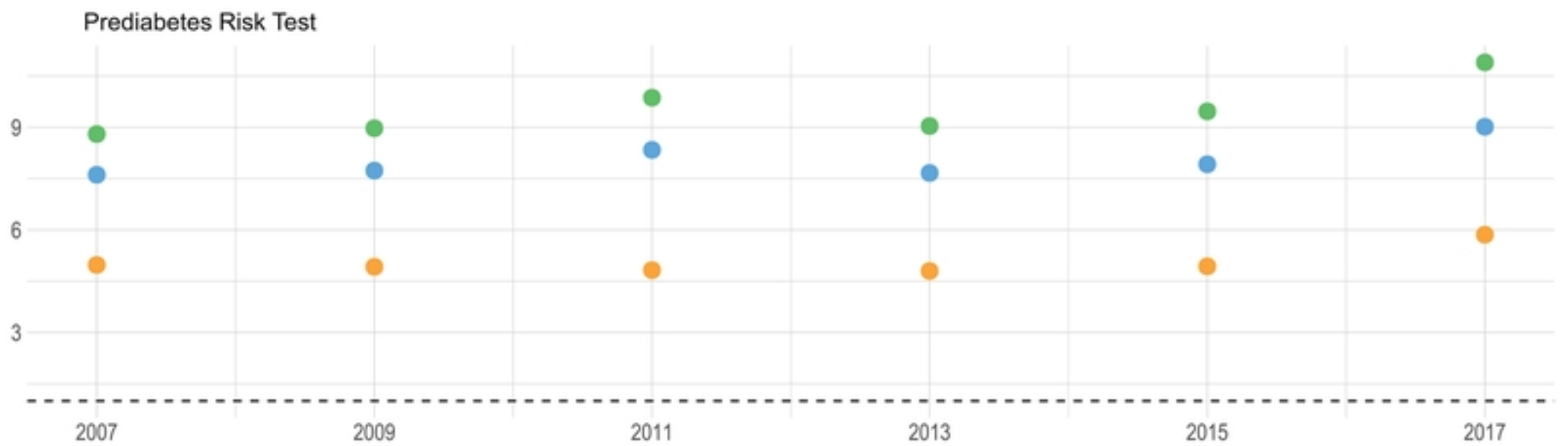
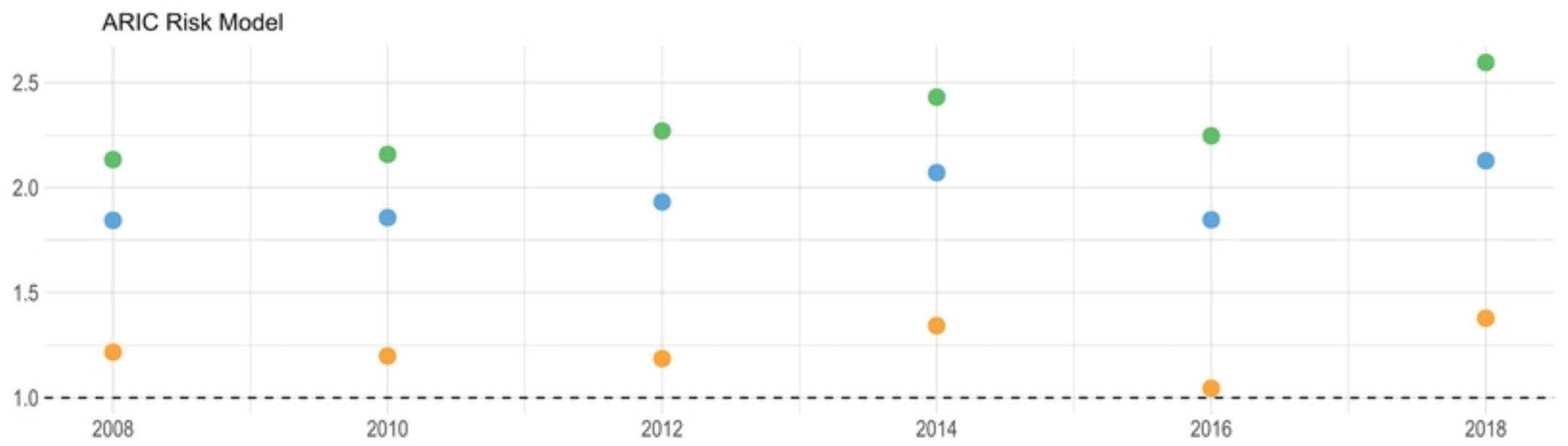
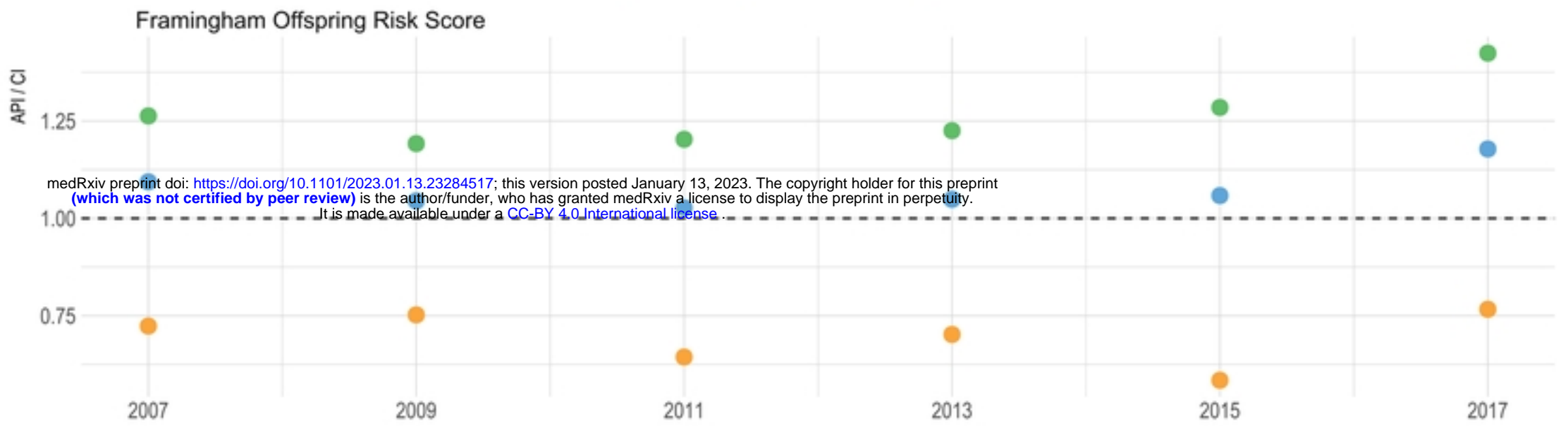
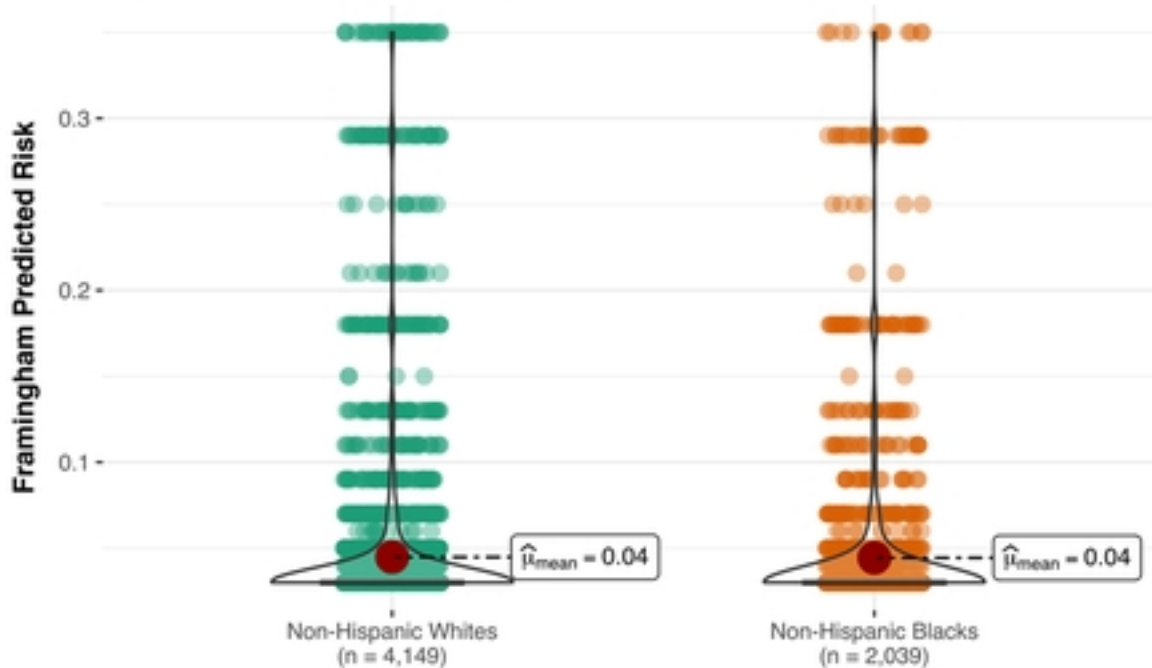


Figure 3

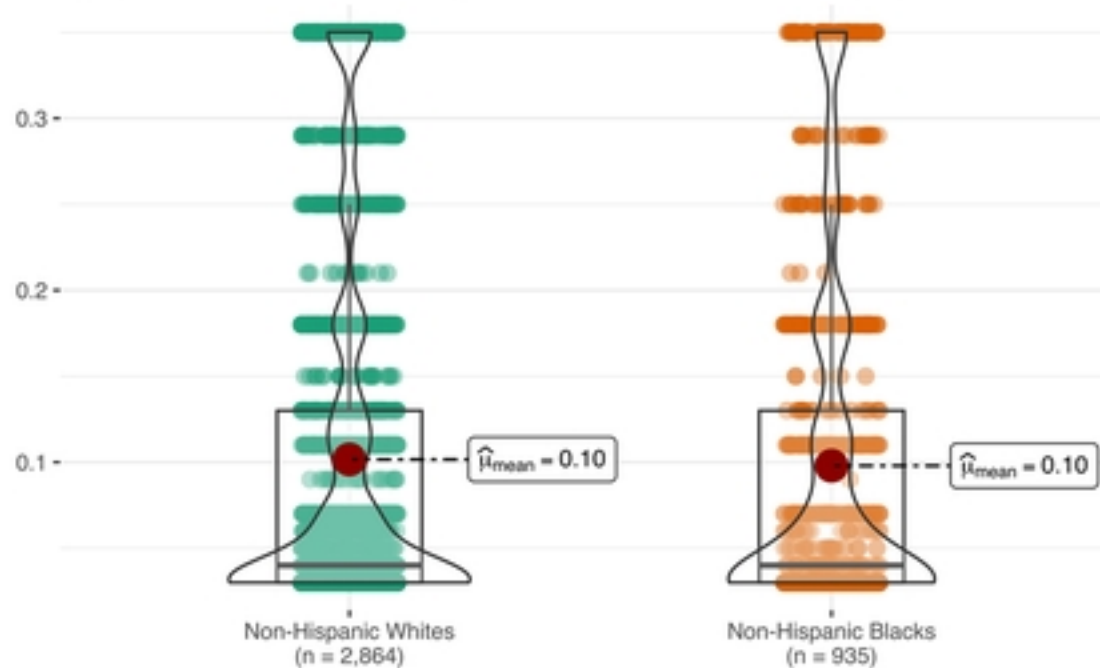
PRT Score < 5

$t_{\text{Welch}}(4211.53) = 0.32, p = 0.75, \hat{g}_{\text{Hedges}} = 8.57\text{e-}03, \text{CI}_{95\%} [-0.04, 0.06], n_{\text{obs}} = 6,188$



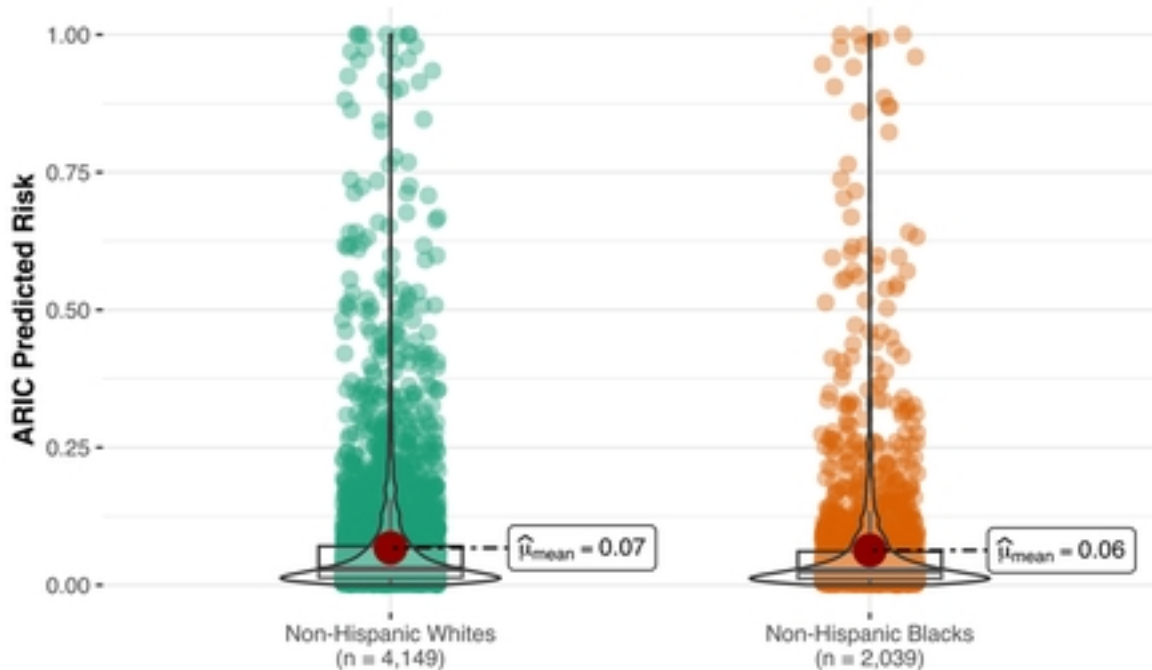
PRT Score ≥ 5

$t_{\text{Welch}}(1640.41) = 0.97, p = 0.33, \hat{g}_{\text{Hedges}} = 0.04, \text{CI}_{95\%} [-0.04, 0.11], n_{\text{obs}} = 3,799$



PRT Score < 5

$t_{\text{Welch}}(3937.63) = 1.38, p = 0.17, \hat{g}_{\text{Hedges}} = 0.04, \text{CI}_{95\%} [-0.02, 0.09], n_{\text{obs}} = 6,188$



PRT Score ≥ 5

$t_{\text{Welch}}(1545.95) = 1.33, p = 0.18, \hat{g}_{\text{Hedges}} = 0.05, \text{CI}_{95\%} [-0.02, 0.12], n_{\text{obs}} = 3,799$

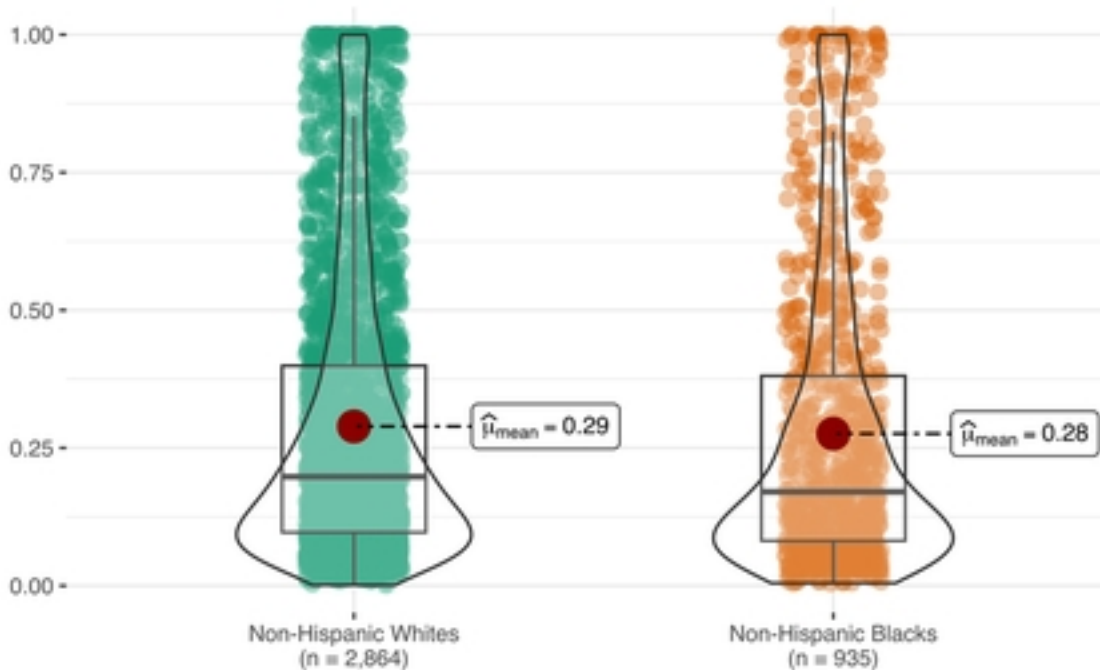


Figure 4