

1 **Passive Detection of COVID-19 with Wearable Sensors and** 2 **Explainable Machine Learning Algorithms**

3 Matteo Gadaleta¹, Jennifer M. Radin¹, Katie Baca-Motes¹, Edward Ramos^{1,2}, Vik Kheterpal²,
4 Eric J. Topol¹, Steven R. Steinhubl¹, Giorgio Quer¹

5 ¹ Scripps Research Translational Institute, 3344 N Torrey Pines Ct Plaza Level, La Jolla,
6 California 92037 USA

7 ² CareEvolution, 625 N Main St, Ann Arbor, Michigan 48104 USA

8 Correspondence to: Giorgio Quer, Scripps Research Translational Institute, 3344 N Torrey Pines
9 Ct Plaza Level, La Jolla CA 92037 USA; gquer@scripps.edu.

10 **ABSTRACT**

11 Individual smartwatch or fitness band sensor data in the setting of COVID-19 has shown promise
12 to identify symptomatic and pre-symptomatic infection or the need for hospitalization,
13 correlations between peripheral temperature and self-reported fever, and an association between
14 changes in heart-rate-variability and infection. In our study, a total of 38,911 individuals (61%
15 female, 15% over 65) have been enrolled between March 25, 2020 and April 3, 2021, with 1,118
16 reported testing positive and 7,032 negative for COVID-19 by nasopharyngeal PCR swab test.
17 We propose an explainable gradient boosting prediction model based on decision trees for the
18 detection of COVID-19 infection that can adapt to the absence of self-reported symptoms and to
19 the available sensor data, and that can explain the importance of each feature and the post-test-
20 behavior for the individuals. We tested it in a cohort of symptomatic individuals who exhibited
21 an AUC of 0.83 [0.81-0.85], or AUC=0.78 [0.75-0.80] when considering only data before the
22 test date, outperforming state-of-the-art algorithm in these conditions. The analysis of all
23 individuals (including asymptomatic and pre-symptomatic) when self-reported symptoms were
24 excluded provided an AUC of 0.78 [0.76-0.79], or AUC of 0.70 [0.69-0.72] when considering
25 only data before the test date. Extending the use of predictive algorithms for detection of
26 COVID-19 infection based only on passively monitored data from any device, we showed that it
27 is possible to scale up this platform and apply the algorithm in other settings where self-reported
28 symptoms can not be collected.

29

30 INTRODUCTION

31 Frequent monitoring to quickly identify, trace and isolate cases of SARS-CoV-2 is needed to
32 help control the spread of the infection as well as improve individual patient care through the
33 earlier initiation of effective therapies.¹ Frequent diagnostic testing is one important option but
34 suffers from implementation challenges and a lack of accessibility for individuals affected most
35 by COVID-19.² Self-reporting of symptoms has been found to be predictive of a positive test,³
36 and could be used to encourage individuals to get tested earlier. However, such an approach not
37 only requires active engagement of the individual, but also misses the approximately one-third of
38 asymptomatic infected individuals completely, and delays diagnosis in those who are infected
39 but presymptomatic.⁴ On the other hand, passive monitoring is possible with commercial sensor
40 devices measuring biometrics such as resting heart rate,⁵ sleep⁶ or activity, which have been
41 shown to be effective in the detection of COVID-19 versus non-COVID-19 when incorporated in
42 combination with self-reported symptoms.⁷

43 Individual sensor data in the setting of COVID-19 has also shown promise in identifying pre-
44 symptomatic infection,⁸ the need for hospitalization,⁹ correlations between peripheral
45 temperature and self-reported fever,¹⁰ differences in the changes in wearable data between
46 individuals with COVID-19 versus influenza-like-illnesses,¹¹ and an association between
47 changes in heart-rate-variability and infection.¹² These studies focused on a specific device
48 brand, or on a predefined set of signals. However, for a broader use of personal health
49 technologies it is important to design algorithms that are device agnostic and can adapt to the
50 specific data collected by any sensor, including the less costly devices.

51 Our prospective app-based research platform DETECT (Digital Engagement and Tracking for
52 Early Control and Treatment) allows participants to enter self-reported symptoms or COVID-19

53 test results, and to share data from any wearable device that is connected to Google Fit or Apple
54 Health Kit platform. In a previous study, we developed a deterministic algorithm to discriminate
55 between symptomatic individuals testing positive or negative for COVID-19, analyzing changes
56 in daily values of resting heart rate, length of sleep and amount of activity, together with self-
57 reported symptoms.⁷

58 In order to provide the most accurate early warnings for COVID-19 to all participants for a wide
59 variety of wearable devices, we proposed and validated a machine learning algorithm that ingests
60 all available sensor data for the detection of COVID-19 infection. The algorithm can outperform
61 our previously proposed algorithm in similar conditions (AUC=0.83, IQR=[0.81, 0.85]), and
62 more importantly, it can automatically adapt to the specific sensor used, exploiting all the
63 information collected from the more advanced sensors or focusing on a smaller set of signals
64 from more basic sensors, and explaining the feature importance and the post-test behavioral
65 changes for the individual. The algorithm uses self-reported symptoms when they are available,
66 or otherwise makes its inference based on sensor data only, thus adapting to different
67 engagement levels of the individuals in the study.

68 **RESULTS**

69 In this study, we investigated the accuracy of a machine learning model in the detection of
70 COVID-19 infection based on the available data acquired from wearable devices and self-
71 reported surveys. We analyzed the accuracy of the detection algorithm for individuals who self-
72 reported at least one symptom prior to the COVID-19 test (named the “symptomatic cohort” in
73 what follows) and not-reporting any symptom prior to the COVID-19 test (the “no-symptom-
74 reported cohort”). We also separately investigated the accuracy obtained using only data
75 collected before a COVID-19 test versus including pre- and post-test data, in order to explore the
76 effect of behavioral changes just the act of testing for COVID-19 might have on individuals.

77 **Participant characteristics**

78 A total of 38,911 individuals (61% female, 15% over 65) have been enrolled between March 25,
79 2020 and April 3, 2021. Among these participants, 1,118 (66% female, 8% over 65) reported at
80 least one positive and 7,032 (63% female, 14% over 65) at least one negative COVID-19 nasal
81 swab test. The total number of COVID-19 swab tests reported during the same period was
82 18,175, with 1,360 (7.5%) positives, 16,398 negatives and 417 with non-reported results. Among
83 the positive tests, 539 (48% of the considered cases) reported at least one symptom in the 15
84 days preceding the test date, 592 (52%) did not report any symptom, and 229 have been excluded
85 from the analysis for lack of sufficient data or for being too close to a prior test.

86 **Dataset description**

87 The participants of the study shared their personal device data (including historical data collected
88 prior to enrollment), self-reported symptoms and diagnostic test results during the data collection
89 period. We divided the measures into four categories: symptom features, including all self-

90 reported symptoms; sensor features, including all measures related to activity, heart rate or sleep;
91 anthropometrics; and demographics. (Table 1)

92 **Detection of COVID-19**

93 The normalized deviations from the baseline for a subset of representative features are reported,
94 (Table 2), highlighting the difference between positive and negative COVID-19 individuals, both
95 excluding and including data after the test date, based on gender and age. As expected, we
96 observed larger variation from the baseline, in terms of heart rate, sleep and activity related
97 features, for individuals who tested positive for COVID-19 with respect to individuals who
98 tested negative. This observation held for all the demographic groups, both excluding or
99 including post-test data. (Figure 1) Based on these features, a prediction model was trained and
100 tested in different conditions.

101 For the symptomatic cohort, we observed a significant difference in the model's output between
102 participants who tested positive or negative, showing that the two groups can be effectively
103 separated, (Figure 1.b) even if we consider only the days preceding the test date (Figure 1.a) thus
104 excluding any behavioral bias potentially caused by taking the test and awaiting results or
105 knowledge of the test outcome. We showed also the predictions for the no-symptom-reported
106 cohort, considering the data before the test date or all the available data, respectively. (Figure 1.c
107 and 1.d) As expected, while a significant difference between the individuals testing positive or
108 negative could still be observed, it is harder to clearly separate the two groups.

109 Symptomatic cases exhibited an area under the receiver operating characteristic (ROC) curve
110 (AUC) of 0.83 [0.81-0.85], while when considering only data before the test date the
111 performance slightly decreased, with AUC=0.78 [0.75-0.80]. For the no-symptom-reported

112 cohort, we observed an AUC of 0.74 [0.72-0.76], or AUC=0.66 [0.64-0.68] when considering
113 only data before the test date. (Figure 2)

114 **Importance of each feature**

115 For the symptomatic cohort, self-reported symptoms were of crucial importance for the most
116 accurate diagnosis of the disease. Considering only data before the test, self-reported symptoms
117 accounted for 60% of the relative contribution to the predictive model, (Figure 3.a) while
118 considering all peri-test data, the importance of the self-reported symptoms decreased to a
119 relative contribution of 46%. (Figure 3.b)

120 For both the symptomatic and no-symptom-reported cohorts, we observed a consistent change in
121 the importance of the activity sensor features, if we consider only data before the test. For the no-
122 symptom-reported cohort, (Figure 3.c and 3.d) the importance of the activity sensor features
123 increased from 46% to 54% when all peri-test data were considered – potentially as a
124 consequence of precautionary measures imposed after testing and awaiting results or receiving a
125 positive test outcome. Sleep sensor features importance did not change significantly when post-
126 test data were included for either the cohort reporting symptoms, or those not reporting
127 symptoms, potentially because sleep was less affected by the knowledge of a test result. Sensor
128 features in the heart rate category had a small relative contribution (6%) for the symptomatic
129 cohort, (Figure 3.a) while their contribution increased (18%) in the no-symptom-reported cohort,
130 (Figure 3.c) acquiring more importance in the absence of information about symptoms and when
131 only pre-test data was considered. Anthropometrics, such as height or weight, provided only a
132 small relative contribution, while the contribution of demographic features, such as age or
133 gender, was negligible.

134 Finally, we provided more details about specific symptoms, and how each of them, on average,
135 affects the model's prediction. (Figure 4) We identified highly discriminative symptoms (cough
136 and decrease in taste and smell, with $\geq 10\%$ relative contribution), medium discriminative
137 symptoms (congestion or runny nose, fever, chills or sweating and congestion or runny nose with
138 $<10\%$ and $\geq 5\%$ relative contribution) and low discriminative features (e.g. body aches,
139 headache, fatigue, with $< 5\%$ relative contribution).

140

141 **DISCUSSION**

142 Our machine learning model based on decision trees can discriminate between individuals who
143 tested positive or negative for COVID-19 based on multiple data types collected by wearable
144 devices, demographic information and self-reported symptoms when available. The adaptability
145 of the algorithm to the available data allows us to also study the performance of the algorithm for
146 individuals in the absence of self-reported symptoms, who may account for almost half of
147 COVID-19 positive individuals.⁴ In order to estimate the effects of the behavioral changes due to
148 the act of testing and/or receiving a positive COVID-19 test, we performed a temporal analysis
149 dividing the data collected before and after the date of COVID-19 testing. The model has been
150 shown to perform well for the identification of COVID-19 infection when incorporating data
151 from symptomatic individuals that includes the five days following the date of testing, with an
152 AUC of 0.83 (IQR: 0.81-0.85). By considering only data preceding the test date, we achieved an
153 AUC of 0.78 (IQR: 0.75-0.80) for people who reported symptoms. When available, self-reported
154 symptoms remain the predominant feature category considered by the model in all our test
155 scenarios, demonstrating the importance of an engaging system that allows participants to easily
156 report this information at any time. Among participants with symptoms, we identified cough and
157 decrease in taste and smell as the most highly discriminative symptoms for a COVID-19
158 infection, followed in order of importance by fever, chills or sweating, and congestion or runny
159 nose.

160 Using the same model, we also investigated individuals who did not report any symptoms.
161 Despite the lack of self-reported information about the sickness, the model achieved an AUC of
162 0.74 (IQR: 0.72-0.76) when considering the period following the test, and an AUC of 0.66 (IQR:
163 0.64-0.68) excluding post-test data. Looking at the importance of the features used by the

164 algorithm, we noticed that the importance of sensor-based Activity substantially increases when
165 considering also post-test data, likely reflecting a potential behavioral change for the participants
166 due to imposed precautionary measures.¹³ On the other hand, the importance of the heart rate
167 features, which are less likely to be affected by short term behavioral changes, is higher when the
168 model consider only data before the COVID-19 test. Moreover, since heart rate elevation might
169 serve as an indicator of inflammatory conditions,¹⁴ its relative importance increases significantly
170 in the absence of self-reported symptoms.

171 These results build on our prior retrospective work on resting heart rate⁵ and sleep,⁶ which when
172 aggregated at the population level, have been shown to significantly improve real-time
173 predictions for influenza-like illness.¹⁵ In an early study, using the initial data from DETECT, we
174 demonstrated the potential of using self-reported symptoms and wearable data for the
175 discrimination of positive and negative cases of COVID-19,⁷ which has been validated by
176 several subsequent independent studies evaluating detection of COVID-19 from wearable
177 devices.¹⁶ The availability of high frequency intra-day data has shown promise to identify pre-
178 symptomatic infection⁸, even if additional studies with a larger number of individuals are needed
179 to prove this point. Several studies focused on the specific data provided by a single sensor
180 brand, showing that an increase in respiratory rate¹⁷ and heart rate,⁹ or a decrease in heart rate
181 variability,¹² are significant during an illness, and that the changes in these physiological signals
182 are more severe for COVID-19 positive cases relative to those affected by other influenza-like-
183 illnesses.¹¹

184 We believe a strength of this research program is that anyone, with any wearable sensor, can
185 participate. As wearable sensors continue to evolve and increase in number, predictive
186 algorithms not dependent on a specific device or data type are needed to optimize the value of

187 continuous, individual data. The algorithm proposed in this work is designed to ingest all
188 available data, exploiting the information provided by the most advanced sensors, while
189 detecting the presence of a COVID-19 infection for everybody owning any type of wearable
190 sensor. The algorithm recognized the importance of self-reported symptoms in the prediction
191 accuracy, but it is also designed to work in the absence of them, thus extending its applicability
192 to the asymptomatic, pre-symptomatic or just a less engaged population who may not want to
193 bother with reporting symptoms.

194 The analysis of individuals without self-reported symptoms extends the use of the algorithms for
195 a fully passive monitoring of the pandemic and provides the possibility of applying the algorithm
196 in other settings that collect wearable sensor data but are not equipped to collect and analyze
197 self-reported symptoms. (Supplementary Material) Among them, the largest is Corona-
198 Datspende, a project developed by the Robert Koch Institute to collect sensor data from more
199 than 600,000 individuals, monitoring the course of the pandemic in Germany.¹⁸

200 The negligible importance given by our algorithm to the demographic features may be explained
201 by observing that the physiological features we consider are changes with respect to an
202 individual baseline. While an individual's baseline differs based on their demographic features,
203 the changes with respect to the baseline that we use in our algorithm are not much affected by the
204 demographic characteristics of the individual.

205 While the use of machine learning in the detection and prognostication for COVID-19 based on
206 chest radiographs and CT scans have been questioned in a systematic review that discussed how
207 none of the current studies are of potential clinical use due to biases or methodological flaws,¹⁹
208 the use of machine learning to enable a continuous and passive COVID-19 early detection is
209 both very promising - for the potential to be scaled up effectively to a large fraction of the

210 population - and repeatable - since we used a strictly separated test set for each of the cross-
211 validation folds. Furthermore, the prediction algorithms developed as part of the DETECT
212 system could be adapted to study the long term health problems due to COVID-19,²⁰⁻²⁴ or the
213 effects of COVID-19 vaccine on vital signs and individual behavior.²⁵⁻²⁷ For future infectious
214 pathogen epidemics and pandemics, the new machine learning algorithms developed from the
215 DETECT data can be adapted and re-used for early detection of various types of infections,
216 towards the development of a new system to monitor the spread of future viral illness and
217 prevent future outbreaks or pandemics.

218 **Limitations**

219 In DETECT, all data is participant reported with no validation of the accuracy of self-reported
220 symptoms, test dates or results. While we were able to collect continuous data, the amount of
221 sensor data collected, or the accuracy of self-reported symptoms, depends entirely on the
222 willingness of the participants to wear the sensor and accurately report how they feel. Despite the
223 fact that the information collected may not be as accurate as in a controlled laboratory setting,
224 previous work has demonstrated the value of participant-reported outcomes.²⁸⁻³⁰ In the data
225 analysis, among the people who reported the COVID-19 test outcome (active participants), we
226 separated participants who reported at least one symptom from those who did not report any
227 symptoms. The app indeed did not have an explicit way to report the absence of symptoms, so
228 potentially some symptomatic individual may have not reported their symptoms.

229 Furthermore, this study is based on the aggregation of continuously monitored data into a finite
230 number of daily features. A recent study has provided new insights about the analysis of intra-
231 day changes for monitoring physiological variations,³¹ that may be used in future studies.

232 Changes in more advanced metrics, like respiratory rate,¹⁷ peripheral temperature¹⁰ or HRV,¹²
233 may also prove to add to the prediction of a COVID-19 infection, even if they have been
234 marginally considered in our work since only a small fraction of participants were providing this
235 type of data.

236 While previous studies have shown the importance of remote monitoring of individuals,
237 extending health research beyond the limits of brick and mortar health systems,^{32,33} additional
238 disparities are introduced when the study relies on wearable sensors, due to reduced accuracy for
239 certain skin tones³⁴ and unequal access to this digital technology.³⁵ The decreasing cost of
240 wearable sensors (some now less than \$35) and the inner adaptability of our detection algorithm
241 to any sensor and any given level of engagement of the participant with the in-app system will
242 hopefully help in decreasing the barriers for underserved and underrepresented populations.

243

244

245 **METHODS**

246 **Study Population**

247 Individuals living in the United States and at least 18 years old are eligible to participate in the
248 DETECT study. After downloading the iOS or Android research app, MyDataHelps, and
249 consenting into the study, participants are asked to share their personal device data (including
250 historical data collected prior to enrollment) from any wearable device connected through direct
251 API (for Fitbit devices), or via Apple HealthKit or GoogleFit data aggregators. A participant is
252 invited to report symptoms, diagnostic test results, vaccine status, and connect their electronic
253 health records, but they can opt to share as much or as little data as they would like. The

254 recruitment of participants happens via the study website (www.detectstudy.org), several media
255 reports, or outreach from our partners at Walgreens, CVS/Aetna, Fitbit and others.

256 **Ethical Considerations**

257 All individuals participating in the study provided informed consent electronically. The protocol
258 for this study was reviewed and approved by the Scripps Office for the Protection of Research
259 Subjects (IRB 20–7531).

260 **Data collection, aggregation, and group definition**

261 All the participants with at least one self-reported result for a COVID-19 swab test during the
262 entire data collection period have been considered in this study. Based on the reported data, an
263 individual is considered Negative if the test resulted negative and no other positive tests have
264 been reported in the period from 60 days before to 60 days after the test date. A minimum
265 distance of 60 days is guaranteed between tests from the same individual considered in the
266 analysis. This ensures that, if multiple tests are reported in the same period, only the first one is
267 considered in our analysis, and the ones reported in the following 60 days will be ignored.

268 For each participant, we collect the data preceding and following the test date from all the
269 connected devices, including, among others, detailed sleep intervals, number of steps and daily
270 resting heart rate values. All the considered metrics are reported and detailed. (Table 1) If
271 multiple values per day are available for the same data type, a specific pre-processing has been
272 applied to obtain a single representative daily value. Data has been collected from all the devices
273 synchronized with the Fitbit or HealthKit application available on the smartphone. If data of the
274 same type is available from multiple devices, only the most used device in the monitored period
275 is considered.

276 Along with device data, we also analyze the reported surveys looking for self-reported
277 symptoms. We considered all the symptoms reported from 15 days before to the day of test,
278 further dividing the participants into two groups: Symptomatic cohort, if we observe at least one
279 reported symptom before the day of test, and non-symptom reporting cohort, if no symptom has
280 been reported during this period. The frequency of each reported symptom for positive and
281 negative cases are also reported. (Figure 5)

282 **Baseline evaluation**

283 Behavioral and physiological data acquired from wearable devices are highly subjective. The
284 intrinsic inter-individual variability of physiological metrics, the different habits of the users, and
285 the multiple purposes of the wearable devices requires a careful definition of the subjective
286 baseline value for each of the considered metrics. Thus, the daily baseline is calculated using an
287 exponentially weighted moving average:

$$Baseline[d] = \sum_{n=0}^{60} Weight[n] \times DailyValue[d - n]$$

288 where d is the current day and n represent the number of days before d , with a maximum of 60
289 days before the current date, while $DailyValue$ can be any of the daily data measures among the
290 ones considered.

291 The oscillation of the measures during the baseline period is also subjective to change over time.
292 To measure the daily baseline variability, we evaluate the weighted standard deviation using the
293 same weights of the baseline

$$BaselineVariability[d] = \sqrt{\sum_{n=0}^{60} Weight[n] \times (DailyValue[d - n] - Baseline[d])^2}$$

294 The weights (Figure 6) decrease exponentially as $e^{-\alpha n}$ with $\alpha=0.05$. We exclude the first 7 days
295 ($n < 7$) from the computation to avoid recent changes to affect the baseline value.
296 Many behavioral habits present strong weekly patterns, such as an increased sleep duration
297 during the weekend, or weekly physical activities. To take into account these behaviors in the
298 baseline evaluation, and to reduce the chances of false positives, we consider weekly patterns by
299 increasing three times all the weights corresponding to the same day of the week.
300 During the course of a temporary disease, physiological measures may be different. In order not
301 to affect the baseline value, which should only depend on the normal behavior, we exclude the
302 10 days following any reported symptoms from the baseline evaluation.
303 Finally, the weights are normalized to sum to 1.

$$w[n] = \begin{cases} 0 & \text{if } n < 7 \text{ or } n > 60 \\ e^{\alpha n} & \text{if } n > 7 \\ 3e^{\alpha n} & \text{if } n = 7m, \quad m = 1, 2, \dots, 8 \\ 0 & \text{if } n \in [s, s + 10], \quad s = \text{symptom date} \end{cases}$$
$$Weight[n] = \frac{w[n]}{\sum_i w[i]}$$

304 The deviation from the baseline values is then evaluated as:

$$DailyValueDev[n] = \frac{DailyValue[n] - Baseline[n]}{BaselineVariability[n]}$$

305 This value represents how far the specific metric is from the expected normal value, day by day.
306 Values are considered to be valid only if at least 50% of corresponding data are available during
307 the baseline period.

308 **Feature extraction**

309 We propose two analyses, one considering all available data (5 days before and 5 days after the
310 test date), and one considering only the period before the test date (5 days before test), in order to

311 further analyze the impact of the test outcome on the individual behavior and the natural course
312 of the disease.

313 We consider four different macro categories of feature. (Table 1)

314 - *Sensor features*: all the features acquired or derived from the device measurements
315 belong to this group. In this study, we consider the minimum, average, and maximum
316 deviation values from the baseline in the days considered. This category is further divided
317 into 3 sub-categories, including activity, heart and sleep related features.

318 - *Symptom features*: a separate binary feature is considered for each of the reported
319 symptoms. If the corresponding symptom has been reported in the considered period its
320 value is set to 1, otherwise 0.

321 - *Anthropometrics*: if available from the monitored devices, several anthropometric
322 features are also considered like body weight, height, body mass index, fat percentage,
323 and basal metabolic rate.

324 - *Demographic features*: this category includes age and gender self-reported by the
325 participants.

326 Using the aforementioned features, we developed a gradient boosting prediction model based on
327 decision trees³⁶. The model has been trained and tested in four different conditions, using data
328 from the symptomatic or no-symptom-reported cohort, and preceding the reported test date or
329 considering all available data around the test date. Normalized deviation (Z-score) from a
330 subjective and dynamic baseline value was evaluated daily for each metric and each individual.
331 A weighted average based on past data was defined as the baseline estimation, whose weights are
332 reported. (Figure 6)

333 The entire dataset has been randomly divided into 5 separate non-overlapping test sets. For each
334 test set, a model is trained using all the remaining data, ensuring an equal percentage of positive
335 cases between train and test sets. For each model, we also ensure that the test set remains strictly
336 separate from the training, so training data are not involved in the test.

337 To analyze the intrinsic variability of the model due to data availability, we estimate 95%
338 confidence intervals for the presented results. Bootstrap method has been utilized for this
339 purpose, with 10,000 independent random iterations from the test outcomes.

340 To have a better understanding of the effect of COVID-19 on physiological and behavioral
341 aspects, we consider symptomatic and no-symptom-reported cases separately. Additionally,
342 different models have been analyzed considering sensor features evaluated including and
343 excluding sensor data after the reported test date. Comparative results are presented in terms of
344 AUC of the ROCs. Sensitivity (SE), specificity (SP), positive predictive value (PPV) and
345 negative predictive value (NPV), associated to an optimal operating point, are also reported. The
346 optimal operating point is defined as the point with the highest average value between SE and
347 SP.

348 The interpretable nature of the decision tree model allows for the evaluation of feature
349 importance estimates^{37,38}. To this end, we evaluate, for each feature, the average prediction
350 changes when the feature value is perturbed. The higher the change to the prediction value, the
351 higher the contribution given by the corresponding feature to the model's outcome. To have a
352 more comprehensive overview of the feature importance, we further aggregated the importance
353 associated to features in the same category.

354

355 **ACKNOWLEDGEMENTS**

356 This work was funded by grant number UL1TR002550 from the National Center for Advancing
357 Translational Sciences (NCATS) at the National Institutes of Health (NIH) (E.J.T., S.R.S., G. Q.)
358 and by grant number OIA-2040727 (Convergence Accelerator) at the National Science
359 Foundation (NSF) (M.G., G.Q.).

360

361 **AUTHORS CONTRIBUTIONS**

362 M.G., and G.Q. made substantial contributions to the study conception and design. M.G., J.M.R.,
363 K.B.-M., E.R., S.R.S, and G.Q. made substantial contributions to the acquisition of data. M.G.,
364 and G.Q. conducted statistical analysis. M.G., S.R.S., and G.Q. made substantial contributions to
365 the interpretation of data. M.G., and G.Q. drafted the first version of the manuscript. M.G.,
366 J.M.R., K.B.-M., E.R., V. K., E.J.T., S.R.S., and G.Q. contributed to critical revisions and
367 approved the final version of the manuscript. M.G., and G.Q. take responsibility for the integrity
368 of the work.

369

370 **COMPETING INTERESTS**

371 S.R.S. is employed by PhysIQ. The other authors declare no competing interests.

372

373

Feature Category	Feature Description	Total COVID-19 Individual Cases	Number of Individuals	Fitbit Users [%]	Available Days Median [IQR]
Symptoms Features	Fatigue	1149	1091	-	-
	Headache	1104	1061	-	-
	DifficultyBreathing	208	206	-	-
	DiarrheaOrVomiting	378	368	-	-
	DecreaseInTasteSmell	247	247	-	-
	Cough	892	859	-	-
	FeverChillsOrSweating	660	645	-	-
	CongestionOrRunnyNose	1152	1097	-	-
	NeckPain	409	396	-	-
	BodyAches	823	795	-	-
	SoreThroat	973	922	-	-
StomachAche	312	302	-	-	
Sensor Features (Activity)	Total number of daily steps	9348	6983	82 %	673 [507 - 707]
	Total daily distance traveled on foot	9281	6938	82 %	676 [554 - 708]
	Calories burned from periods above sedentary level	7629	5679	100 %	676 [554 - 708]
	Calories burned inclusive of BMR	8279	6142	100 %	676 [554 - 708]
	Minutes spent fairly active	7291	5450	100 %	676 [554 - 708]
	Minutes spent sedentary	8196	6082	100 %	676 [554 - 708]
	Minutes spent very active	7171	5370	100 %	676 [554 - 708]
	Minutes spent lightly active	7628	5678	100 %	676 [554 - 708]
Sensor Features (Heart)	Daily resting heart rate	9105	6810	81 %	492 [341 - 653]
	Maximum daily heart rate variability	1812	1407	0 %	363 [185 - 499]
	Minimum daily heart rate variability	1812	1407	0 %	363 [185 - 499]
	Average daily heart rate variability	1812	1407	0 %	363 [185 - 499]
Sensor Features (Sleep)	Total daily sleep time	7473	5603	94 %	445 [240 - 629]
	Total daily time spent in bed	7473	5603	94 %	445 [240 - 629]
	Sleep efficiency of the main sleep	7473	5603	94 %	445 [240 - 629]
	Sleep time of the main sleep	7473	5603	94 %	445 [240 - 629]
Anthropometrics	Body Mass Index	8478	6303	-	-
	Self-reported height	1673	1277	-	-
	Body weight	9240	6896	-	-
	Body fat percentage	4594	3431	-	-
	Basal metabolic rate (BMR) only calories	8279	6142	-	-
Demographic Features	Self-reported gender	10494	7853	-	-
	Age at the time of test	10479	7841	-	-

374

375 *Table 1 - Description and categorization of the features. In the table, one case refers to a specific test*
 376 *from an individual, while an individual may report multiple tests. The total number of individual*
 377 *COVID-19 cases is the number of cases with the corresponding feature value available for the analysis.*
 378 *The available days for Sensor Features represent the median number of days available for all the*
 379 *participants (IQR reported in brackets).*

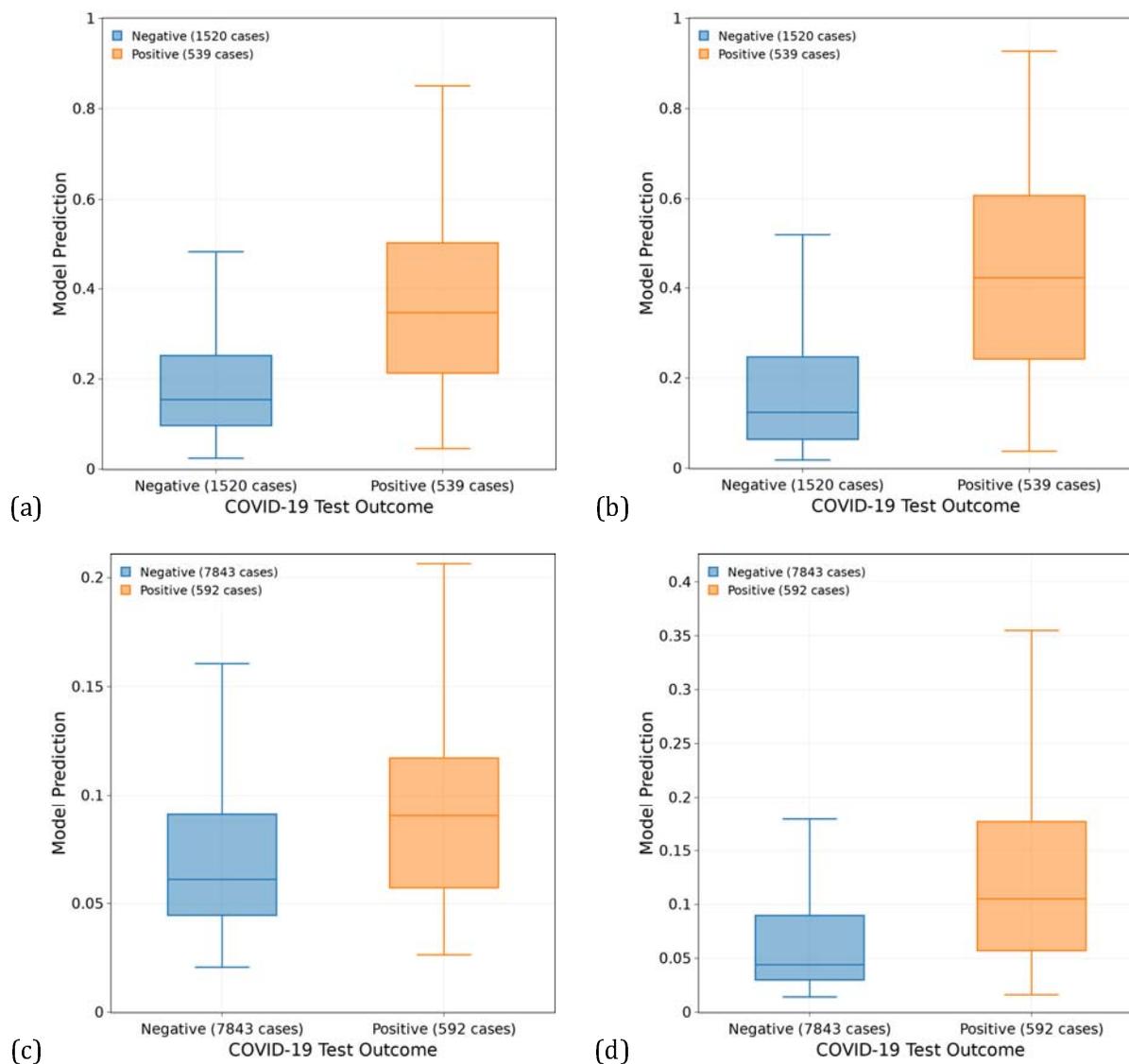
380

Feature	Subset	Z-score: Median [IQR] (number of valid samples)			
		Excluding data after COVID-19 Test		Including data after COVID-19 Test	
		COVID-19 Negative	COVID-19 Positive	COVID-19 Negative	COVID-19 Positive
Steps	Overall	-0.08 [-0.45, 0.34] (8290)	-0.44 [-0.84, 0.01] (977)	-0.10 [-0.44, 0.27] (8360)	-0.68 [-1.10, -0.24] (988)
	Male	-0.08 [-0.44, 0.34] (2977)	-0.42 [-0.84, -0.01] (324)	-0.10 [-0.43, 0.26] (3004)	-0.65 [-1.02, -0.28] (326)
	Female	-0.07 [-0.46, 0.34] (5313)	-0.46 [-0.85, 0.03] (653)	-0.10 [-0.44, 0.27] (5356)	-0.70 [-1.13, -0.22] (662)
	Young (< 40)	-0.05 [-0.44, 0.35] (2759)	-0.39 [-0.76, 0.17] (348)	-0.06 [-0.40, 0.31] (2793)	-0.65 [-1.00, -0.10] (354)
	Middle Age (40-65)	-0.08 [-0.45, 0.33] (4543)	-0.47 [-0.88, -0.04] (557)	-0.11 [-0.45, 0.24] (4575)	-0.72 [-1.15, -0.31] (561)
	Old Age (> 65)	-0.10 [-0.47, 0.35] (975)	-0.41 [-0.80, 0.04] (71)	-0.14 [-0.50, 0.26] (979)	-0.51 [-1.14, -0.01] (72)
Resting Heart Rate	Overall	0.04 [-0.38, 0.49] (8061)	0.23 [-0.26, 0.78] (945)	0.04 [-0.33, 0.44] (8150)	0.16 [-0.28, 0.70] (955)
	Male	0.05 [-0.33, 0.48] (2922)	0.27 [-0.21, 0.77] (318)	0.05 [-0.29, 0.44] (2955)	0.20 [-0.20, 0.75] (321)
	Female	0.04 [-0.42, 0.50] (5139)	0.22 [-0.27, 0.77] (627)	0.04 [-0.36, 0.44] (5195)	0.12 [-0.32, 0.65] (634)
	Young (< 40)	0.06 [-0.44, 0.56] (2685)	0.26 [-0.29, 0.79] (333)	0.06 [-0.38, 0.49] (2725)	0.11 [-0.36, 0.68] (337)
	Middle Age (40-65)	0.04 [-0.36, 0.48] (4403)	0.23 [-0.26, 0.76] (542)	0.04 [-0.31, 0.42] (4446)	0.17 [-0.25, 0.69] (548)
	Old Age (> 65)	0.01 [-0.33, 0.36] (960)	0.17 [-0.15, 0.73] (69)	0.02 [-0.28, 0.36] (966)	0.10 [-0.09, 0.79] (69)
Average Daily Heart Rate Variability	Overall	-0.02 [-0.41, 0.39] (1591)	-0.25 [-0.65, 0.19] (196)	0.02 [-0.33, 0.32] (1614)	-0.11 [-0.46, 0.29] (198)
	Male	-0.02 [-0.38, 0.38] (743)	-0.27 [-0.68, 0.14] (84)	0.02 [-0.30, 0.31] (757)	-0.13 [-0.61, 0.18] (84)
	Female	-0.00 [-0.43, 0.40] (848)	-0.21 [-0.60, 0.20] (112)	-0.00 [-0.36, 0.34] (857)	-0.10 [-0.43, 0.37] (114)
	Young (< 40)	-0.04 [-0.43, 0.45] (469)	-0.20 [-0.60, 0.39] (71)	0.01 [-0.36, 0.34] (473)	-0.04 [-0.41, 0.43] (72)
	Middle Age (40-65)	-0.01 [-0.38, 0.37] (885)	-0.29 [-0.67, 0.13] (110)	0.01 [-0.32, 0.32] (901)	-0.10 [-0.46, 0.17] (111)
	Old Age (> 65)	0.01 [-0.39, 0.40] (235)	-0.19 [-0.39, -0.01] (15)	0.02 [-0.32, 0.28] (238)	-0.20 [-0.51, -0.03] (15)
Daily Sleep Time	Overall	0.01 [-0.36, 0.36] (6575)	0.23 [-0.20, 0.71] (791)	0.00 [-0.28, 0.30] (6672)	0.38 [-0.01, 0.85] (801)
	Male	0.00 [-0.36, 0.35] (2299)	0.25 [-0.20, 0.71] (255)	-0.00 [-0.29, 0.29] (2340)	0.34 [-0.03, 0.83] (258)
	Female	0.01 [-0.36, 0.37] (4276)	0.23 [-0.19, 0.70] (536)	0.00 [-0.28, 0.30] (4332)	0.40 [-0.00, 0.86] (543)
	Young (< 40)	0.02 [-0.36, 0.40] (2259)	0.21 [-0.17, 0.66] (270)	0.01 [-0.27, 0.32] (2294)	0.39 [0.03, 0.84] (276)
	Middle Age (40-65)	0.00 [-0.36, 0.35] (3555)	0.26 [-0.20, 0.73] (463)	0.00 [-0.28, 0.30] (3606)	0.39 [-0.03, 0.88] (467)
	Old Age (> 65)	-0.04 [-0.39, 0.29] (750)	0.11 [-0.17, 0.67] (57)	-0.06 [-0.31, 0.25] (761)	0.30 [-0.02, 0.66] (57)
		BMI: Median [IQR] (number of valid samples)			
		COVID-19 Negative		COVID-19 Positive	
BMI	Overall	26.56 [23.52, 31.00] (7561)		27.60 [23.92, 32.03] (917)	
	Male	26.65 [24.11, 30.28] (2679)		27.80 [24.50, 30.57] (303)	
	Female	26.50 [23.13, 31.51] (4882)		27.46 [23.48, 32.46] (614)	
	Young (< 40)	26.07 [22.97, 30.93] (2635)		26.93 [23.45, 31.74] (346)	
	Middle Age (40-65)	27.18 [23.96, 31.43] (4075)		27.99 [24.38, 32.29] (501)	
	Old Age (> 65)	25.64 [23.14, 28.91] (840)		27.53 [23.46, 30.29] (69)	

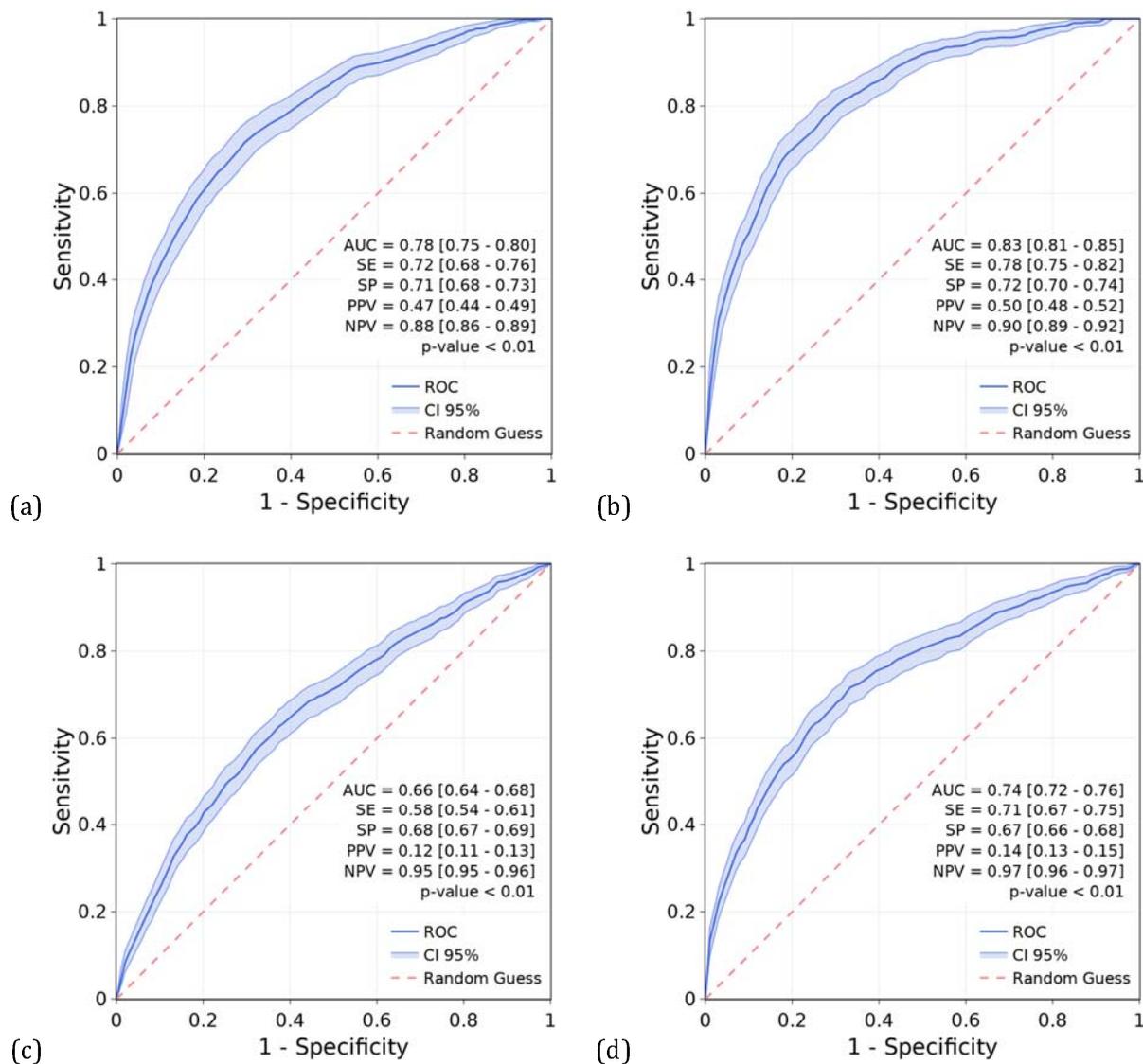
381
382 *Table 2 – Normalized deviation from baseline values for a selected number of representative features.*

383 *Values for positive and negative COVID-19 individuals, including or excluding the period after the test*
384 *date, are reported. The results are stratified among gender and age groups. Median, interquartile*
385 *range (IQR) and number of COVID-19 cases analyzed are reported.*

386



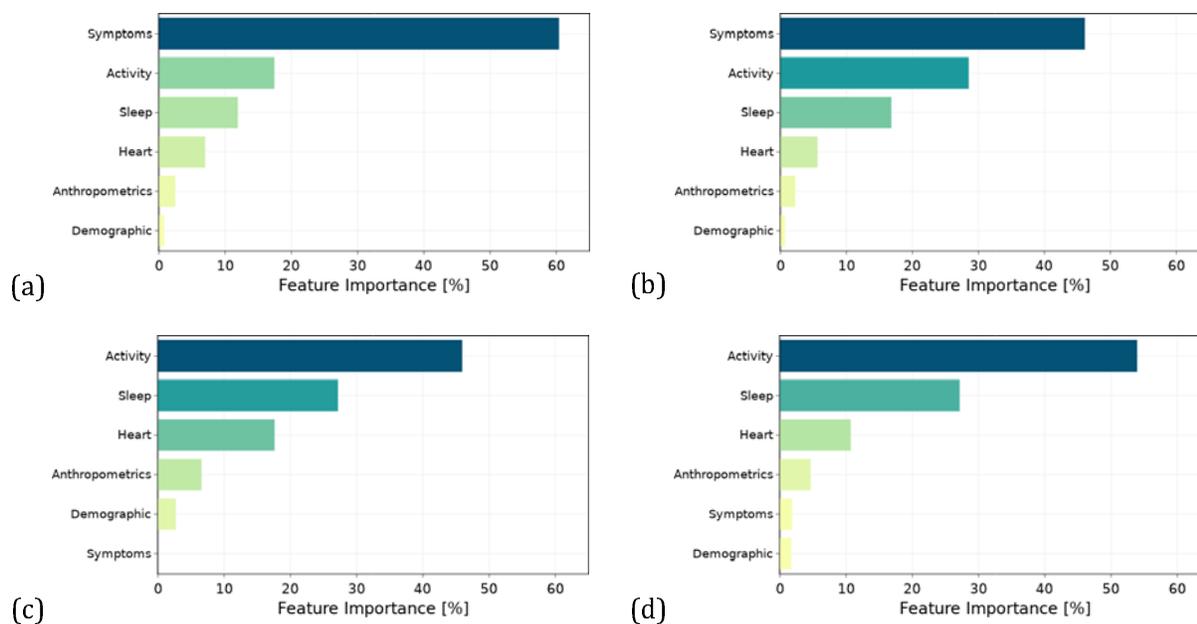
387 *Figure 1 – Output of the prediction models for symptomatic cases, excluding (a) and including the data*
388 *after the test date (b), and for no-symptom-reported cases, excluding (c) and including the data after*
389 *the test date (d). The boxes represent the IQR, and the horizontal lines are the median values. The*
390 *number of cases considered for the analysis are reported in the legend.*



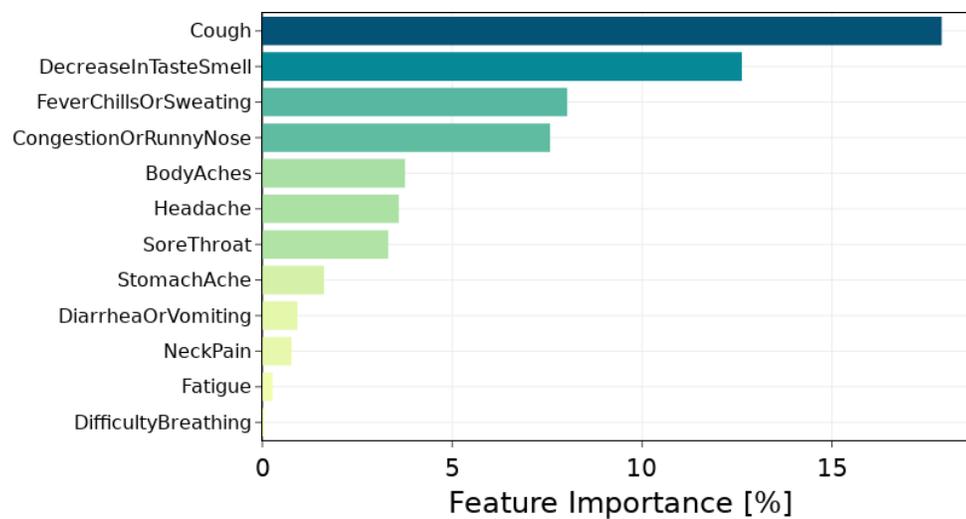
391 *Figure 2 - Receiver operating characteristic curves (ROCs) for the discrimination between COVID-19*
392 *positive and COVID-19 negative. Performance for symptomatic cases, excluding (a) and including the*
393 *data after the test date (b), and for no-symptom-reported cases, excluding (c) and including the data*
394 *after the test date (d), are reported. The model is a gradient boosting prediction model based on*
395 *decision trees. Median values and 95% confidence intervals (CIs) for sensitivity (SE), specificity (SP),*
396 *positive predictive value (PPV) and negative predictive value (NPV) are reported, considering the*

397 *point on the ROC with the highest average value of sensitivity and specificity. Error bars represent*

398 *95% CIs. P-values of the one-sided Mann-Whitney U test are reported.*



399 *Figure 3 – Overall feature importance based on the average prediction changes when the feature value*
400 *is perturbed. Values are normalized as percentages. Features have been aggregated into macro*
401 *categories. Results for symptomatic cases, excluding (a) and including data after test date (b), and for*
402 *no-symptom-reported cases, excluding (c) and including data after test date (d), are reported.*



403

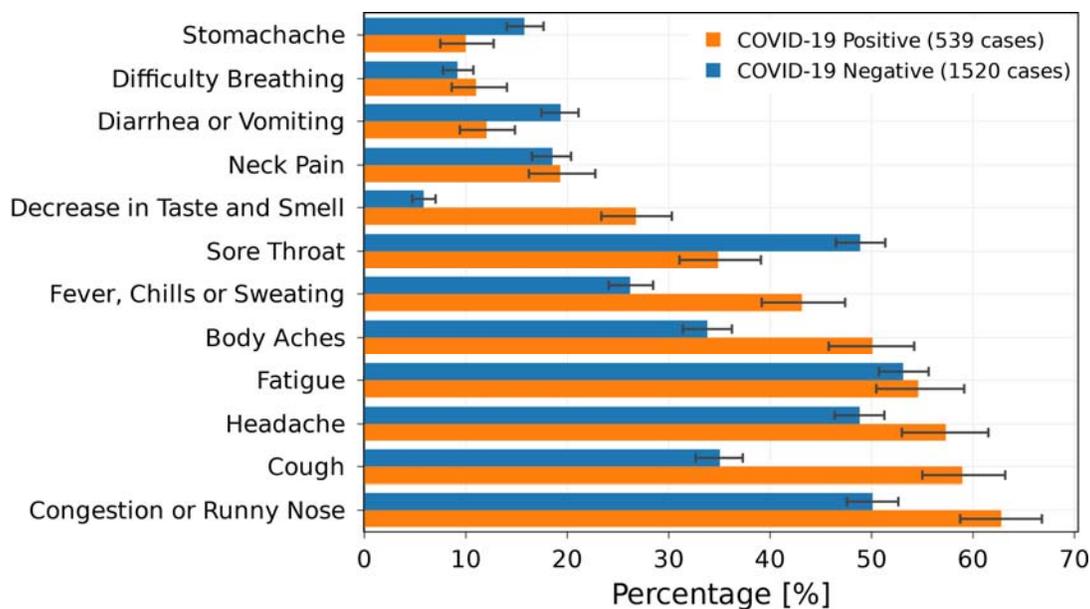
404 *Figure 4 - Feature importance associated to specific symptoms. Only symptoms reported before the*
405 *test date have been considered. Values are normalized as percentages. The results refer to*
406 *symptomatic cases only.*

407

408

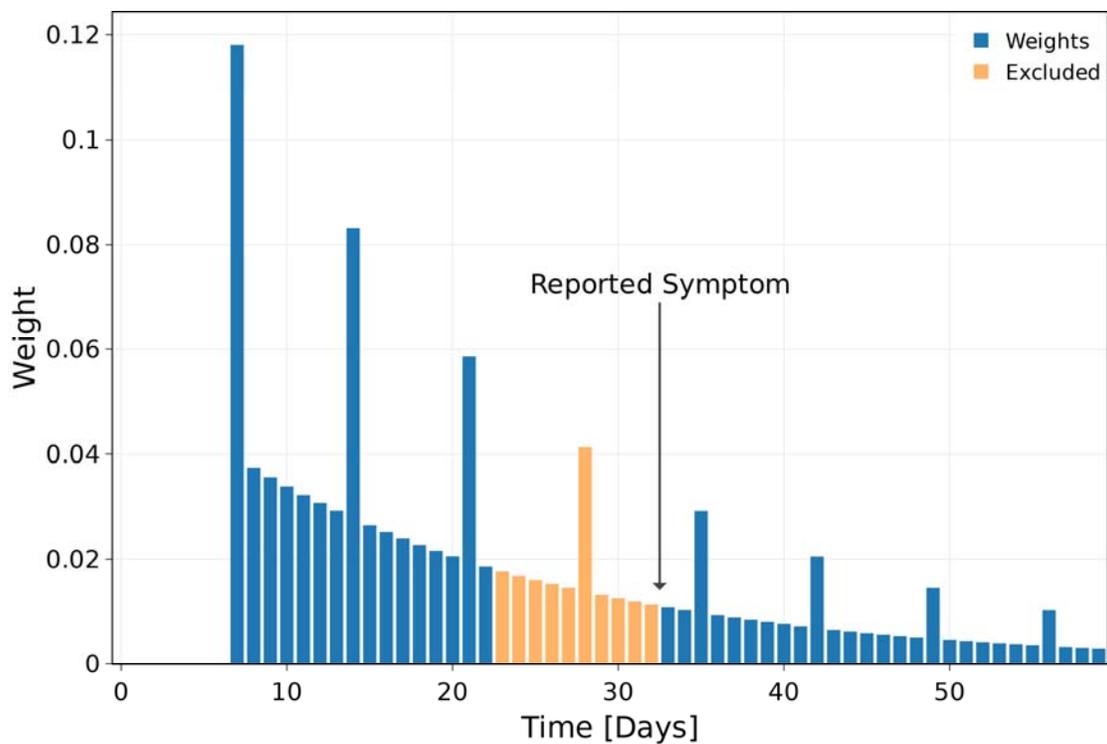
409

410



411

412 *Figure 5 – Percentage of reported symptoms for participants who reported at least one symptom from*
413 *15 days before to 5 days after the test date. The frequencies of the indicated symptoms are shown for*
414 *positive and negative cases. The error bars represent 95% percent confidence intervals.*



415

416 *Figure 6 – Exponentially decreasing weights for the evaluation of the baseline data, with weekly*
417 *patterns. The abscissa represents the temporal distance preceding the analyzed day considered for the*
418 *baseline evaluation. The first 6 days have been excluded to avoid recent changes to affect the baseline.*
419 *Additionally, if a symptom has been reported in this time frame, we set the weights to zero from the*
420 *day of symptom to the next 10 days.*

421 **REFERENCES**

- 422 1. NIH. COVID-19 Treatment Guidelines.
423 <https://www.covid19treatmentguidelines.nih.gov/whats-new/>. (2021).
- 424 2. Manabe, Y.C., Sharfstein, J.S. & Armstrong, K. The Need for More and Better Testing
425 for COVID-19. *Jama* **324**, 2153-2154 (2020).
- 426 3. Menni, C., *et al.* Real-time tracking of self-reported symptoms to predict potential
427 COVID-19. *Nature Medicine* (2020).
- 428 4. Oran, D.P. & Topol, E.J. Prevalence of Asymptomatic SARS-CoV-2 Infection. *Annals of*
429 *Internal Medicine* (2020).
- 430 5. Quer, G., Gouda, P., Galarnyk, M., Topol, E.J. & Steinhubl, S.R. Inter- and
431 intraindividual variability in daily resting heart rate and its associations with age,
432 sex, sleep, BMI, and time of year: Retrospective, longitudinal cohort study of 92,457
433 adults. *PloS one* **15**, e0227709 (2020).
- 434 6. Jaiswal, S.J., *et al.* Association of Sleep Duration and Variability With Body Mass
435 Index: Sleep Measurements in a Large US Population of Wearable Sensor Users.
436 *JAMA Internal Medicine* (2020).
- 437 7. Quer, G., *et al.* Wearable sensor data and self-reported symptoms for COVID-19
438 detection. *Nature Medicine* **27**, 73-77 (2021).
- 439 8. Mishra, T., *et al.* Pre-symptomatic detection of COVID-19 from smartwatch data.
440 *Nature Biomedical Engineering* (2020).
- 441 9. Natarajan, A., Su, H.-W. & Heneghan, C. Assessment of physiological signs associated
442 with COVID-19 measured using wearable devices. *npj Digital Medicine* **3**, 156
443 (2020).
- 444 10. Smarr, B.L., *et al.* Feasibility of continuous fever monitoring using wearable devices.
445 *Sci Rep* **10**, 21640 (2020).
- 446 11. Shapiro, A., *et al.* Characterizing COVID-19 and Influenza Illnesses in the Real World
447 via Person-Generated Health Data. *Patterns* **2**, 100188 (2021).
- 448 12. Hirten, R.P., *et al.* Physiological Data from a Wearable Device Identifies SARS-CoV-2
449 Infection and Symptoms and Predicts COVID-19 Diagnosis: Observational Study. *J*
450 *Med Internet Res* (2021).
- 451 13. Cleary, J.L., Fang, Y., Sen, S. & Wu, Z. A Caveat to Using Wearable Sensor Data for
452 COVID-19 Detection: The Role of Behavioral Change after Receipt of Test Results.
453 *medRxiv*, 2021.2004.2017.21255513 (2021).
- 454 14. Whelton, S.P., *et al.* Association between resting heart rate and inflammatory
455 biomarkers (high-sensitivity C-reactive protein, interleukin-6, and fibrinogen) (from
456 the Multi-Ethnic Study of Atherosclerosis). *The American journal of cardiology* **113**,
457 644-649 (2014).
- 458 15. Radin, J.M., Wineinger, N.E., Topol, E.J. & Steinhubl, S.R. Harnessing wearable device
459 data to improve state-level real-time surveillance of influenza-like illness in the USA:
460 a population-based study. *The Lancet Digital Health* **2**, e85-e93 (2020).
- 461 16. Radin, J.M., Quer, G., Jalili, M., Hamideh, D. & Steinhubl, S.R. The hopes and hazards of
462 using personal health technologies in the diagnosis and prognosis of infections. *The*
463 *Lancet Digital Health* (to appear, 2021).

- 464 17. Miller, D.J., *et al.* Analyzing changes in respiratory rate to predict the risk of COVID-
465 19 infection. *PloS one* **15**, e0243693 (2020).
- 466 18. Robert Koch-Institut. Corona Datenspende, [https://corona-](https://corona-datenspende.de/science/en)
467 [datenspende.de/science/en](https://corona-datenspende.de/science/en). (2020).
- 468 19. Roberts, M., *et al.* Common pitfalls and recommendations for using machine learning
469 to detect and prognosticate for COVID-19 using chest radiographs and CT scans.
470 *Nature Machine Intelligence* **3**, 199-217 (2021).
- 471 20. Dani, M., *et al.* Autonomic dysfunction in 'long COVID': rationale, physiology and
472 management strategies. *Clin Med (Lond)* **21**, e63-e67 (2021).
- 473 21. Puntmann, V.O., *et al.* Outcomes of Cardiovascular Magnetic Resonance Imaging in
474 Patients Recently Recovered From Coronavirus Disease 2019 (COVID-19). *JAMA*
475 *Cardiology* **5**, 1265-1273 (2020).
- 476 22. Sudre, C.H., *et al.* Attributes and predictors of long COVID. *Nature Medicine* (2021).
- 477 23. Logue, J.K., *et al.* Sequelae in Adults at 6 Months After COVID-19 Infection. *JAMA*
478 *Network Open* **4**, e210830-e210830 (2021).
- 479 24. Radin, J.M., *et al.* Assessment of Prolonged Physiological and Behavioral Changes
480 Associated with COVID-19 Infection. *JAMA Network Open* (to appear, 2021).
- 481 25. Voysey, M., *et al.* Safety and efficacy of the ChAdOx1 nCoV-19 vaccine (AZD1222)
482 against SARS-CoV-2: an interim analysis of four randomised controlled trials in
483 Brazil, South Africa, and the UK. *Lancet* **397**, 99-111 (2021).
- 484 26. Benedict, C. & Cedernaes, J. Could a good night's sleep improve COVID-19 vaccine
485 efficacy? *Lancet Respir Med* (2021).
- 486 27. Quer, G., *et al.* The Physiologic Response to COVID-19 Vaccination. *medRxiv*,
487 2021.2005.2003.21256482 (2021).
- 488 28. Basch, E., *et al.* Overall Survival Results of a Trial Assessing Patient-Reported
489 Outcomes for Symptom Monitoring During Routine Cancer Treatment. *Jama* **318**,
490 197-198 (2017).
- 491 29. Bell, S.K., *et al.* Frequency and Types of Patient-Reported Errors in Electronic Health
492 Record Ambulatory Care Notes. *JAMA Network Open* **3**, e205867-e205867 (2020).
- 493 30. Rivera, S.C., *et al.* The impact of patient-reported outcome (PRO) data from clinical
494 trials: a systematic review and critical analysis. *Health and Quality of Life Outcomes*
495 **17**, 156 (2019).
- 496 31. Mishra, T., *et al.* Pre-symptomatic detection of COVID-19 from smartwatch data.
497 *Nature Biomedical Engineering* **4**, 1208-1220 (2020).
- 498 32. Steinhubl, S.R., *et al.* Effect of a Home-Based Wearable Continuous ECG Monitoring
499 Patch on Detection of Undiagnosed Atrial Fibrillation: The mSToPS Randomized
500 Clinical Trial. *Jama* **320**, 146-155 (2018).
- 501 33. Radin, J.M., *et al.* Pregnancy health in POWERMOM participants living in rural versus
502 urban zip codes. *J Clin Transl Sci* **4**, 457-462 (2020).
- 503 34. Colvonen, P.J., DeYoung, P.N., Bosompra, N.-O.A. & Owens, R.L. Limiting racial
504 disparities and bias for wearable devices in health science research. *Sleep* **43**(2020).
- 505 35. Beaunoyer, E., Dupéré, S. & Guitton, M.J. COVID-19 and digital inequalities:
506 Reciprocal impacts and mitigation strategies. *Comput Human Behav* **111**, 106424
507 (2020).

- 508 36. Dorogush, A.V., Ershov, V. & Gulin, A. CatBoost: gradient boosting with categorical
509 features support. *arXiv preprint arXiv:1810.11363* (2018).
- 510 37. Lundberg, S.M. & Lee, S.-I. Consistent feature attribution for tree ensembles. *arXiv*
511 *preprint arXiv:1706.06060* (2017).
- 512 38. Lundberg, S.M., Erion, G.G. & Lee, S.-I. Consistent individualized feature attribution
513 for tree ensembles. *arXiv preprint arXiv:1802.03888* (2018).
- 514
- 515

516 **SUPPLEMENTARY MATERIAL**

517 **Analysis in the absence of self-reported symptoms**

518 For an improved and effective tracking of the pandemic, many research institutions are collecting
519 sensor data from individuals, while it is not always possible to provide surveys and collect active
520 feedback from participants. Information actively added by participants may be crucial especially
521 for potentially infected individuals, indeed a fully passive data collection system may be adopted
522 by a broader audience and have a more capillary diffusion.

523 Our model can be leveraged to support also these studies based uniquely on passive data
524 collection. To this end, we performed an additional analysis without the inclusion of self-
525 reported symptoms as source of knowledge for the model. The results of the analysis without
526 considering self-reported symptoms are shown in terms of AUC of the ROC. (Figure S.1)

527

