

# A Method for Intelligent Allocation of Diagnostic Testing by Leveraging Data from Commercial Wearable Devices: A Case Study on COVID-19

Jessilyn Dunn (✉ [jessilyn.dunn@duke.edu](mailto:jessilyn.dunn@duke.edu))

Duke University <https://orcid.org/0000-0002-3241-8183>

**Md Mobashir Hasan Shandhi**

Duke University

**Peter Cho**

Duke University

**Ali Roghanizad**

Duke University

**Karnika Singh**

Duke University

**Will Wang**

Duke University

**Oana Enache**

Duke University

**Amanda Stern**

Duke University

**Rami Sbahi**

Duke University <https://orcid.org/0000-0002-2081-5894>

**Bilge Tatar**

Duke University

**Sean Fiscus**

Duke University

**Qi Xuan Khoo**

Duke University

**Yvonne Kuo**

Duke University

**Xiao Lu**

Duke University

**Joseph Hsieh**

Duke University

**Alena Kalodzitsa**

Duke University

**Amir Bahmani**

Stanford University <https://orcid.org/0000-0003-4533-9334>

**Arash Alavi**

Stanford University <https://orcid.org/0000-0002-7245-6889>

**Utsab Ray**

Stanford University

**Michael Snyder**

Stanford University School of Medicine <https://orcid.org/0000-0003-0784-7987>

**Geoffrey Ginsburg**

Duke Institute for Genome Sciences & Policy (IGSP) <https://orcid.org/0000-0003-4739-9808>

**Dana Pasquale**

Duke University <https://orcid.org/0000-0001-6686-7844>

**Christopher Woods**

Duke University Medical Center

**Ryan Shaw**

Duke University

---

**Article**

**Keywords:**

**Posted Date:** April 1st, 2022

**DOI:** <https://doi.org/10.21203/rs.3.rs-1490524/v1>

**License:**  This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

---

**Version of Record:** A version of this preprint was published at npj Digital Medicine on September 1st, 2022. See the published version at <https://doi.org/10.1038/s41746-022-00672-z>.

# Abstract

Mass surveillance testing can help control outbreaks of infectious diseases such as COVID-19. However, diagnostic test shortages are prevalent globally and continue to occur in the US with the onset of new COVID-19 variants, demonstrating an unprecedented need for improving our current methods for mass surveillance testing. By targeting surveillance testing towards individuals who are most likely to be infected and, thus, increasing testing positivity rate (i.e., percent positive in the surveillance group), fewer tests are needed to capture the same number of positive cases. Here, we developed an Intelligent Testing Allocation (ITA) method by leveraging data from the CovIdentify study (6,765 participants) and the MyPHD study (8,580 participants), including smartwatch data from 1,265 individuals of whom 126 tested positive for COVID-19. Our rigorous model and parameter search uncovered the optimal time periods and aggregate metrics for monitoring continuous digital biomarkers to increase the positivity rate of COVID-19 diagnostic testing. We found that resting heart rate features distinguished between COVID-19 positive and negative cases earlier in the course of the infection than steps features, as early as ten and five days prior to the diagnostic test, respectively. We also found that including steps features increased the area under the receiver operating characteristic curve (AUC-ROC) by 7–11% when compared with RHR features alone, while including RHR features improved the AUC of the ITA model's precision-recall curve (AUC-PR) by 38–50% when compared with steps features alone. The best AUC-ROC ( $0.73 \pm 0.14$  and  $0.77$  on the cross-validated training set and independent test set, respectively) and AUC-PR ( $0.55 \pm 0.21$  and  $0.24$ ) were achieved by using data from a single device type (Fitbit) with high-resolution (minute-level) data. Finally, we show that ITA generates up to a 6.5-fold increase in the positivity rate in the cross-validated training set and up to a 3-fold increase in the positivity rate in the independent test set, including both symptomatic and asymptomatic (up to 27%) individuals. Our findings suggest that, if deployed on a large scale and without needing self-reported symptoms, the ITA method could improve allocation of diagnostic testing resources and reduce the burden of test shortages.

## Introduction

The COVID-19 pandemic has severely impacted our world, with more than 478 million COVID-19 cases and 6.13 million deaths estimated worldwide [1]. In the US alone, there have been more than 81 million cases and nearly 1 million deaths at the time of writing [2]. Mass surveillance testing has been identified as the most effective tool to monitor the spread of infectious diseases including COVID-19 [3]. However, a combination of cost, availability, and impracticality of frequent and widespread testing impedes the mass epidemiologic surveillance needed to curb new disease outbreaks. To date, COVID-19 diagnostic test shortages are still prevalent globally, and shortages continue to occur in the US with the onset of new variants (e.g, Delta, Omicron) [4]–[6]. For example, when the Delta variant emerged in July 2021, daily demand for tests across the US surged from 250k to 1.5 million in the span of 2 months [7]. A similar circumstance occurred with the Omicron variant, where testing capacity failed to meet the sudden demand [8]–[10]. Furthermore, rural-urban disparities in testing access have exacerbated existing inequities resulting in further harm to underserved communities [11], [12]. In June 2020, it was estimated

that 64% of counties in the United States, predominantly rural, did not have access to COVID-19 testing [13]. Such circumstances lead to underreporting of COVID-19 incidence and may lead to a premature sense of security and unwarranted changes in public health measures [12]. Thus, there is an unprecedented need to improve our current and future methods for mass COVID-19 surveillance testing, especially as stronger testing capacity has been associated with reduced mortality and greater pandemic control [14].

By targeting surveillance testing toward individuals who are more likely to be infected with the disease, more positive cases can be captured with the same number of tests, increasing the positivity rate of the tested population (Fig. 1A) [4]. The positivity rate (i.e., percent positive rate or percent positive) is the percentage of all diagnostic tests performed that are positive. The likelihood of disease presence prior to a diagnostic test, or the pretest probability, is dependent on disease prevalence in the population under surveillance. By filtering the broader surveillance population to a subpopulation with a higher likelihood of infection, the allocation and utility of tests can be improved (Fig. 1A). In other words, more positive cases can be captured with the same number of tests and, thus, the testing positivity rate is increased. The development of tools to increase testing positivity rate are not only crucial in the early phase of a pandemic outbreak when the available clinical diagnostic testing tools are inadequate to meet the existing demand, but also throughout a pandemic in remote locations, underserved communities, and low- and middle-income countries where testing is known to be particularly scarce [15].

The rapid adoption of commercial wearable devices such as smartwatches and activity trackers brings forth opportunities to develop novel tools to support an intelligent disease detection infrastructure. Multiple studies suggest the utility of digital biomarkers, objective and quantifiable digitally-collected physiological and behavioral data (e.g., resting heart rate, step count, sleep duration, and respiratory rate), collected by consumer devices along with patient-reported symptoms to monitor the progression of respiratory and influenza-like illnesses [16]–[22].

To determine who to test in settings where there are a limited number of diagnostic tests available (i.e., limited testing capacity), we explored whether information from wearables could help rank individuals by their likelihood of a current COVID-19 infection. To achieve this, we developed an Intelligent Testing Allocation (ITA) model which leverages differences in digital biomarkers to distinguish individuals who are likely positive or negative for COVID-19 in order to improve current methods of diagnostic test allocation and increase testing positivity rates.

## Results

We developed the CovIdentify platform in April 2020 to integrate commercial wearable device data and electronic symptom surveys to calculate an individual's real-time risk of being infected with COVID-19. A total of 7,348 participants e-consented to the CovIdentify study between April 2, 2020 and May 25, 2021 through the secure research electronic data capture (REDCap) system (Fig 1B) [23]. Of those consented, 6,765 participants enrolled in the study (Supplementary Table 1) by completing an enrollment survey

consisting of 37-61 questions that followed branching logic (Supplementary Document 1) [23]. Of those enrolled, 2,887 participants connected their smartwatches to the CovIdentify platform, including 1,689 Garmin, 1,091 Fitbit, and 107 Apple smartwatches. Throughout the course of the study, 362,108 daily surveys were completed by 5,859 unique participants, with a mean of 62 and a median of 37 daily surveys completed per individual. Of all CovIdentify participants, 1,289 participants reported at least one diagnostic test result for COVID-19 (132 positive and 1,157 negative) (Fig 1B). All survey and device data collected through CovIdentify was transferred securely to a protected cloud environment for further analysis. Out of the 1,289 participants with self-reported diagnostic test results, 136 participants (16 positive and 120 negative) had smartwatch data available during the time periods needed for analysis. The relatively small number of participants with available smartwatch data out of the larger population is consistent with similar bring-your-own-device (BYOD) studies aimed toward COVID-19 infection prediction from personal devices [18], [19], [24].

***Development of the Intelligent Testing Allocation (ITA) model.*** A diagnostic testing decision support model was designed to leverage real-world data to intelligently allocate diagnostic tests in a surveillance population where there are insufficient tests available to test all people in the surveillance group (Fig 1A, top). To increase the study population size, we augmented our dataset with data from the MyPHD study. Similar to CovIdentify, MyPHD collected simultaneous smartwatch, symptom, and diagnostic testing data during the COVID-19 pandemic [24], [25]. The wearables and diagnostic testing data were publicly available [24], [25] while symptom data were added for this work. From the MyPHD study, smartwatch, symptom, and diagnostic testing data from an additional 1,129 participants (110 positive and 1,019 negative) were included in this analysis.

***Differences in resting heart rate (RHR) and steps measured by smartwatches well before and immediately prior to a COVID-19 diagnostic test.***

To compare digital biomarkers between healthy and infected states, data were segmented into two time periods: a baseline period (22 to 60 days prior to the diagnostic test date) and a detection period (21 days prior to the diagnostic test date). We chose this window for the detection period to encompass the COVID-19 incubation period (2-14 days) reported by the CDC as well as the common delay between symptom onset and diagnostic testing. Consistent with prior literature [16], [26], daily RHR increased significantly during the detection period from baseline for those who were COVID-19 positive, with an average difference ( $\pm$  SD) of  $1.65 \pm 4.63$  bpm ( $n=117$ ,  $p$ -value  $< 0.001$ , paired t-test) over the entire time periods. On average, daily RHR values more than two standard deviations from the baseline mean were present as early as 13 days prior to the positive test, with an increasing trend that peaked at one day prior to the test date (Fig 1C, bottom). Conversely, the step count during the detection period decreased significantly from baseline, with a difference of  $-854 \pm 2386$  steps/day ( $n=125$ ,  $p$ -value  $< 0.0001$ ). On average, step counts less than two standard deviations from the baseline mean were present as early as 10 days prior to the

positive test and reached the minimum value 2 days after the test date (Fig 1C, top). For the subset of participants in our dataset with available symptom onset dates, daily RHR and step count that differed beyond two standard deviations from the baseline mean occurred as early as five days before the symptom onset date (Supplementary Fig 1). Timelines for this and other real-world infection studies should be considered as rough estimates because exact dates of exposure and symptom onset are unknown, unlike in controlled infection studies [22], [27]. Our findings, however, are consistent with the 2-14 day COVID-19 incubation period reported by the CDC [28].

There was also a significant difference in digital biomarkers between the baseline and detection periods of participants who tested negative, but it was less pronounced than for those who tested positive. Specifically, the daily RHR difference was  $0.58 \pm 4.78$  bpm ( $n=1,094$ ,  $p$ -value  $< 0.05$ ) and the step count difference was  $-281 \pm 2,013$  steps/day ( $n=1,136$ ,  $p$ -value  $< 0.0001$ ). We hypothesized that the digital biomarker differences in the COVID-19 negative group were because a subset of the negative group may have experienced a health anomaly other than COVID-19 (e.g., influenza) that resulted in physiological differences between the baseline and detection periods. Another recent study also observed RHR elevation and activity reduction in individuals who were COVID-19 negative but flu positive, and the magnitudes of these differences were lower than in individuals who were COVID-19 positive [18]. To explore the possibility that our COVID-19 negative group contains false negatives due to test inaccuracies or physiological differences due to a health anomaly besides COVID-19, we performed hierarchical clustering on the symptom data from individuals who reported negative tests and found a trend toward multiple subgroups (Extended Data Fig 1). This finding supports the existence of COVID-19 negative subgroups. It should also be noted that the highly significant  $p$ -value for the digital biomarker differences in the COVID-19 negative group is likely attributable to the higher number of participants (9-fold higher) compared with the COVID-19 positive group.

**Cohort Definition.** For the ITA model development, we only included subjects with sufficient wearable data (>50% data availability) in each of the baseline and detection periods. Consequently, 83 participants from CovIdentify (9 COVID-19 positive and 74 COVID-19 negative) and 437 participants from MyPHD (54 COVID-19 positive and 383 COVID-19 negative) were included in the ITA model development process (Table 1). Of the 520 total participants with sufficient wearable data, 469 participants had high frequency minute-level wearable data (280 from Fitbits) from which we calculated daily RHR and step counts. Device-reported daily values were available for the remaining 51 participants. To explore whether high frequency wearable data, or high frequency wearable data from a single device type, could improve the performance of digital biomarkers for ITA, we developed and validated our ITA model using three cohorts, which we refer to as 1) the All-Frequency (AF) cohort: participants with both high frequency and device reported daily values, 2) the All-High-Frequency (AHF) cohort: participants with high frequency data only, and 3) the Fitbit-High-Frequency (FHF) cohort: participants with high frequency Fitbit data only (Extended

Data Fig 2). We analyzed these three cohorts separately in the subsequent analysis and compared the resulting ITA model performance. We divided each cohort into an 80% train and 20% test split, with FHF as a subgroup of AHF, which itself is a subset of AF to ensure that no observations in the training set of one cohort existed in the test set of another (Extended Data Fig 2).

*Feature engineering:* To explore differences in digital biomarkers (median or mean) between the detection and baseline periods that may be useful for the development of ITA model features, we designed four deviation metrics including: (1)  $\Delta$  (detection – baseline), (2) normalized  $\Delta$ , (3) standardized  $\Delta$ , and (4) z-score ((detection - baseline mean) / baseline standard deviation) (Table 2). Each of the four deviation metrics were calculated on the training data by digital biomarkers (RHR and step count), day in the detection period, and cohort (examples in Supplementary Fig 2 and 3), resulting in four calculated metrics per cohort per biomarker. These training data deviation metrics were used as inputs into the subsequent statistical analysis for feature extraction and the ITA model training. We extracted the same resultant features from the independent test set for subsequent ITA model evaluation.

On average, step count decreased ( $\Delta$ Steps) significantly from baseline to the detection period in COVID-19 positive versus negative participants (574 vs. 179, 479 vs. 234, and 601 vs. 216 steps per day for the AF, AHF, and FHF training data, respectively; p-value < 0.05, unpaired t-tests) (Fig 2A, Extended Data Fig 3A, and Extended Data Fig 4A top plots). Conversely, RHR increased ( $\Delta$ RHR) significantly from baseline to the detection period in COVID-19 positive versus negative participants (1.8 vs. 0.7, 1.9 vs. 0.8, and 1.8 vs. 0.7 bpm for the AF, AHF, and FHF training data, respectively; p-value < 0.05, unpaired t-test) (Fig 2A, Extended Data Fig 3A, and Extended Data Fig 4A bottom plots). The 95% confidence intervals of the mean  $\Delta$ Steps and the mean  $\Delta$ RHR overlap considerably between positive and negative participants for the initial phase of the detection period (approximately twenty-one to five days prior to the test date). However, closer to the diagnostic test date (approximately four to one days prior to the test date) the 95% confidence intervals of mean  $\Delta$ Steps largely do not overlap, and the 95% confidence intervals of mean  $\Delta$ RHR do not overlap at all (Fig 2A). The fact that the 95% confidence intervals of mean  $\Delta$ Steps and mean  $\Delta$ RHR do not overlap later in the detection period is consistent with prior literature [29] and suggests that it is possible to aggregate data into summary statistics to develop a decision boundary that effectively separates COVID-19 positive and negative cases. However, the overlap in estimated mean values prior to day 5 suggests that separation between positive and negative cases may be more challenging prior to that point in time. Although the 95% confidence intervals closer to the test date were non-overlapping, there was overlap in the variance of the digital biomarkers between the two groups during that time period (Extended Data Figure 5), which may hinder model performance as separation of the 95% confidence intervals does not necessarily imply significant differences between the groups [30]. Similar estimates of variability have not been reported prior, so we were unable to compare our mean statistics variability to prior literature.

*Optimizing the detection period for the ITA model:* To maximize the separability of the COVID-19 positive and negative groups in the training set, we performed statistical analysis to explore how different lengths and start times of the detection window, parametrized respectively by two variables (the detection end date (DED), defined by days prior to the diagnostic test date, and the detection window length (DWL) defined by number of days), would affect the separation between these two groups. We performed a combinatorial analysis across these two parameters (DED and DWL) to calculate five summary statistics (mean, median, maximum, minimum, and range) of the four deviation metrics (Table 2) to be used as features for model building. This resulted in 40 total summary statistics (20 each from steps and RHR), which we refer to as steps and RHR features, respectively. Statistical comparison of the steps and RHR features between the COVID-19 positive and COVID-19 negative groups was performed on the training data for the AF, AHF, and FHF cohorts separately to uncover the statistically significant features (unpaired t-tests; Benjamini-Hochberg corrected p-value < 0.05).

A systematic grid search to optimize the detection end date and detection window length (DED and DWL, respectively) demonstrated that the closer the detection period is to the diagnostic test date, the larger the number of features that are significantly different between the COVID-19 positive and negative groups (Fig 2B, Extended Data Fig 3B and Extended Data Fig 4B). Across all DED values, DED of one day prior to the diagnostic test date (DED = -1) generated the largest number of significant features for all cohorts. Also, across all cohorts, there were more significant RHR features than steps features (Fig 2B, Extended Data Fig 3B and 4B). Additionally, RHR features became significant earlier in the detection period than steps features (DED as early as -10 days vs. -5 days, respectively), which indicates that changes in RHR occur earlier than steps during the course of infection. Comparison across the three cohorts revealed AF generated the highest number of significant features compared with the AHF and FHF cohorts, which may be attributable to the larger population size of AF. This demonstrates the tradeoff in wearables studies between high frequency data, which is less common but contains more information, and larger population data, which contains data at a variety of sampling frequencies but overall more data to train the models. Across the DWL values, 3 and 5 days generated the largest number of significant features for all cohorts (Figure 2C, Extended Data Fig 3C and 3C), while 5 days also corresponded to the date of the maximum divergence between  $\Delta$ Steps and  $\Delta$ RHR (Fig 2A). Ultimately, this systematic analysis pointed to an optimal DED of 1 day prior to the diagnostic test date and an optimal DWL of 5 days for the detection window duration, both of which were used to generate features for the ITA model.

*ITA feature selection:* When implementing the DED timepoint and DWL duration that best separated the COVID-19 positive and negative groups, there were 28-31 significant features (p-value < 0.05; unpaired t-tests with Benjamini-Hochberg multiple hypothesis correction) that overlapped across the three cohorts, indicating their robustness to differences in data resolution and device types (Supplementary Table 2).

The top 7-9 features, ranked in order of significance, originated from the RHR digital biomarker. To gain a more mechanistic understanding of the RHR and step digital biomarkers, we explored the top two most significantly different (lowest p-value) features for each digital biomarker between those who were COVID-19 positive or negative in the AF cohort (Fig 2D). The decrease in steps during the detection period as compared to baseline was greater in those with COVID-19, with a 2054 vs 99 median decrease in steps (median  $\Delta$ Steps) and a 1775 vs 64 mean decrease in steps for those who were COVID-19 positive vs those who were COVID-19 negative, respectively (p-values < 0.0001). Conversely, the increase in maximum deviation in RHR in the detection period as compared to baseline (maximum  $\Delta$ RHR) and the increase in mean of Z-scores in the detection period as compared to baseline (mean of Z-score RHR) were both significantly higher for COVID-19 positive participants compared to COVID-19 negative participants (8.4 vs. 4.3 bpm for maximum  $\Delta$ RHR and 0.9 vs. 0.2 for the mean of Z-score-RHR; p-values < 0.0001). Consistent across all three cohorts, the median and mean  $\Delta$ Steps were the most significant (lowest p-value) steps features (Extended Data Fig 3D and 4D). However, the top two RHR features differed, which were median and mean Z-score-RHR, and maximum  $\Delta$ RHR and maximum of normalized  $\Delta$ RHR for the AHF and FHF cohorts, respectively (Extended Data Fig 3D and 4D and Supplementary Table 2). The observation of the same top two steps features given the differences in the top two RHR features across the three cohorts may originate from the resolution and device-reported digital biomarkers. For example, the definition of a step and the calculation of the daily step count may be more similar across different device types, while the RHR definition and available HR data resolution may vary more substantially across device types. Although these top features are significantly different between those who are COVID-19 positive and negative, their distributions do overlap, even though the tailedness varies in direction and extent (Fig 2D, Extended Data Fig 3D and 4D, and Supplementary Fig 4), which points to broader challenges surrounding predictive modeling efforts using standard consumer wearable device data for COVID-19 infection detection.

*Development of the ITA model:* To achieve our broader goal of determining who should receive a diagnostic test under circumstances where there are limited tests available, we aimed to design a model that outputs the probability of a person being infected. However, because our ground truth information is binary (positive or negative for COVID-19), we designed this model as a binary classifier which enabled straightforward evaluation of its performance. We used the features that were significantly different in the training data between those who were COVID-19 positive and negative (29 features for AF, 28 for AHF, and 31 for FHF) as inputs into five machine learning classification models: logistic regression, k-nearest neighbors (KNN), support vector machine (SVM), Random Forest (RF), and extreme gradient boosting (XGBoost) (Extended Data Table 1). We chose these five well-established classification models to explore how increasing model complexity and the addition of non-linearity impacts the model performance. We trained these classification models on the training data using nested cross-validation with an inner loop for hyperparameter tuning and an outer loop for model selection. We chose recall as our preferred scoring metric for model selection and evaluation to emphasize the relative impact/cost of false negatives

compared to false positives, as an individual who is truly positive for COVID-19 and is wrongly classified as negative (or healthy) would further spread disease.

Following training, we evaluated the performance of the trained model on the independent test set and used two well-established reporting metrics, including the most commonly reported metric for studies of this kind (the area under the curve for the receiver operating characteristic curve (AUC-ROC)) [26], [31]–[35], and the metric that is most appropriate for this classification task (AUC for the precision-recall curve (AUC-PR)) [36] (Extended Data Table 1; Fig 3, Fig 4, and Extended Data Fig 6). AUC-PR is more appropriate with class-imbalanced data [36], [37], which is the case here (12-15% COVID-19 positive and 85-88% negative in each of the three cohorts). The class imbalance in our dataset was not correctable through resampling methods - we have observed that distributions of features overlap between the COVID-19 positive and negative participants, as demonstrated in the individual feature comparison (Fig 2D, Extended Data Fig 3D and 4D), as well as in the low dimensional representation (using principal component analysis and t-stochastic neighbor embedding) of all the features in the training set of the AF cohort (Supplementary Fig 5).

Of the five models tested, logistic regression outperformed all other models based on the training AUC-PR for all three cohorts and was also the best performing model based on the training AUC-ROC for the AF and FHF cohorts. The superior performance of the logistic regression among other (more complex and nonlinear) models may be attributed to the tendency of more complex and non-linear models to overfit on the training data [38], which comes to light with our cross-validation methods. The superior performance of the logistic regression also points to the potential to develop explainable machine learning predictive models for the ITA model which enables rapid translation from bench to bedside. Overall, the classifier performed best on the FHF cohort (Extended Data Table 1, Fig 4 A, B, C, and D), followed by the AHF cohort, (Extended Data Fig 6 A, B, C, and D) and finally the AF cohort (Fig 3 A, B, C, and D). These performance differences indicate that device-related and data resolution differences may confound disease-related physiological differences captured by digital biomarkers. Therefore, building models using a single device type and with higher resolution data improves performance. For the FHF cohort, the logistic regression model resulted in an AUC-ROC of  $0.73 \pm 0.12$  and AUC-PR of  $0.55 \pm 0.21$  on the cross-validated training set, and AUC-ROC of 0.77 and AUC-PR of 0.24 on the test set (Fig 4). The AUC-ROC from the models were similar to those reported in recent similar studies [26], [32], [35].

However, the performance of the models based only on AUC-ROC in the context of imbalanced data can be misleading, as a large change in the number of false positives may have a small effect on the false positive rate [37]. The precision metric, which integrates both true positives and false positives, can mitigate the effect of an imbalanced dataset (e.g., the higher proportion of negatives seen in this type of

data) on a model's performance. Our precision-recall analysis (Fig 3B, 4B, and Extended Data Fig 6B) demonstrates that we can improve the recall (minimizing false negatives) at the expense of precision. In an extreme example, we were able to achieve 100% recall with a precision of 0.4 on the cross-validated training set of the FHF cohort, whereas, a dummy classifier with random chance (i.e., Random Testing Allocation (RTA)) can achieve a precision of 0.15 on this dataset. It is also important to note that we are not considering resource-limited settings in the ROC and PR analysis; instead, it is assumed that there are a sufficient number of diagnostic tests available for the entire surveillance group. In a resource-limited setting, 100% recall may not be achievable due to the shortage of diagnostic testing.

To understand the relative contribution of the steps and RHR digital biomarkers to the ITA model performance, we developed two separate sets of models using features based only on either steps or RHR using the training set data with logistic regression. Consistent with previous literature [26], [32] the models using steps-based features alone had a higher AUC-ROC than models using RHR-based features alone (cross-validated AUC-ROC of 0.67 vs. 0.64, 0.69 vs. 0.63, and 0.72 vs. 0.68 for steps vs RHR features for the AF, AHF, and FHF training sets, respectively) (Extended Data Fig 7). Interestingly, when using the AUC-PR as the performance metric, models using features based on RHR digital biomarkers outperformed models using features based on steps digital biomarkers, a finding which has not been previously reported (cross-validated AUC-PR of 0.30 vs. 0.38, 0.28 vs. 0.37, and 0.40 vs. 0.49 for steps and RHR features for the AF, AHF, and FHF training datasets, respectively) (Extended Data Fig 7). Overall, the addition of steps features increased the AUC-ROC of the ITA model by 7-11% compared with RHR features alone, while RHR features improved the AUC-PR of the ITA model by 38-50% compared with steps features alone. These results suggest that, while steps features provide more salient information on the trade-off between the true positive rate and false positive rate, RHR features provide more salient information on the trade-off between the true positive rate and the precision (positive predictive value). In other words, while steps features improved the specificity of the predictive model, RHR features improved the precision. We also compared the relative feature importances in the logistic regression and found that two, one, and four of the top five features originated from RHR in the AF, AHF, and FHF cohorts, respectively, with the remaining features originating from steps (Extended Data Fig 8). In all three cohorts, median  $\Delta$ Steps and mean  $\Delta$ Steps were the two most important steps features, which was consistent with our earlier statistical analysis. Maximum  $\Delta$ RHR was the most important RHR feature for the AF and AHF cohorts and the second most important RHR feature for the FHF cohort, and was also one of the top two most significant features in our earlier statistical analysis for the AF and FHF cohorts.

***Improvement in positivity rate for COVID-19 diagnostic testing using the ITA method.*** We next evaluated how the ITA model can improve the current standard of practice for COVID-19 infection surveillance. Under current surveillance testing methods in the US, while some tests are taken due to symptoms or possible exposure, many are taken as precautionary measures for traveling or for surveillance in schools

and workplaces [28]. While such forms of widespread RTA surveillance are beneficial, the positivity rate of widespread diagnostic testing is typically low and, thus, requires sufficient testing capacity in order to prevent testing shortages (e.g., sold out at-home testing kits). Applying an equivalent RTA surveillance approach to our study population results in a 12% positivity rate in both our AF-training (50 COVID-19 positive participants out of 365 participants in total) and AF-test (13 COVID-19 positive participants out of 92 participants in total) datasets. It is important to note that the 12% positivity rate is consistent for all levels of diagnostic testing capacity (0-100% of population). When employing ITA and adding the constraint of limited diagnostic testing capacity (10-30% of population), the testing positivity rate of the cross-validated model increased 2 to 3-fold (21-36% positivity rate) for the training dataset (Fig 3C) and 1.5 to 2.5-fold (19-29% positivity rate) for the testing dataset (Fig 3F).

Comparison of the three cohorts demonstrated that the best performing ITA model stemmed from the FHF cohort and was followed by the AHF cohort (Fig. 4C and 4F, and Extended Data Fig. 6C and 6F). By utilizing ITA and assuming a diagnostic testing capacity at 10-30% of the population, the positivity rate of the FHF and AHF cross-validated training datasets increased by 4-fold (64% positivity rate) and 3-fold (35% positivity rate) when compared to the RTA positivity rates of 15% and 12% for FHF and AHF cohorts, respectively. For the FHF cohort, the positivity rate further increased up to 6.5-fold (100% positivity rate) in the cross-validated training dataset when the diagnostic testing capacity was reduced to 2.5-5% of the population (5-11 diagnostic tests to be allocated to individuals in the training dataset) (Fig 4F). Using the independent test data set, the positivity rate of the FHF and AHF cohorts increased by 1.5 to 3-fold (17-31% positivity rate) and 2 to 3-fold (21-32% positivity rate), respectively, compared to the RTA positivity rate of 11%, when the diagnostic testing capacity was 10-30% of the population. These results indicate the potential of the ITA model to target diagnostic testing resources towards individuals who have a higher likelihood of testing positive (i.e., increasing the positivity rate of diagnostic testing) and enables more efficient allocation of testing capacity.

We further explored how the ITA model performs in symptomatic versus asymptomatic COVID-19 positive individuals in each cohort. We considered participants to be symptomatic who reported any symptoms in the detection period or on the diagnostic test date. Assuming a diagnostic testing capacity of 30%, ITA indicates testing for 19 of 29 symptomatic and 7 of 21 asymptomatic COVID-19 positive individuals in the cross-validated model, and 5 of 8 symptomatic and 1 of 5 asymptomatic COVID-19 positive individuals in the independent test set of the AF cohort. In other words, 7 of 26 (27%) and 1 of 6 (17%) COVID-19 positive individuals were asymptomatic in the ITA determined subpopulation for the cross-validated training set and independent test set of the AF cohort, respectively. Results were similar for the AHF and FHF cohorts (Extended Data Table 2). These findings indicate that the ITA model can not only target diagnostic testing resources towards individuals with symptoms, but also to those without any reported symptoms, further increasing the utility of this method.

## Discussion

The COVID-19 pandemic revealed the fragility of our existing healthcare infrastructure to detect the virus and prevent its spread. One key tool for reducing disease spread is bringing diagnostic testing to the right people at the right time and ensuring appropriate interpretation of the diagnostic testing results based on the prevalence of the disease in the population [4]. In light of this need, in April 2020 we developed CovIdentify to integrate commercial wearable device data and electronic symptom surveys to assess the real-time risk of being infected with COVID-19. We envisioned two possible scenarios where CovIdentify would be useful for informing intelligent testing decisions, including 1) ranking individuals in a group by likelihood of current infection with COVID-19 to determine *who* to test, and 2) tracking a single individual over time for evidence of new infection onset to determine *when* to test. In our initial development of the Intelligent Testing Allocation (ITA) model, we focused on the first question, and ultimately improved the positivity rate of COVID-19 diagnostic testing up to 6.5-fold when compared against Random Testing Allocation (RTA). Based on these results, if deployed on a large scale, the ITA model could be used to better allocate diagnostic testing resources. This method is likely applicable to other diagnostic areas as well, where digital biomarkers can be used to indicate the likelihood of disease.

In this work, we demonstrated that wearable device data can be used to strategically target the allocation of diagnostic tests to where they are most useful. This approach not only increases testing efficiency and allocation but also reduces the costs and supply chain burden of surveillance testing which is an ongoing challenge. Our results further demonstrate that the ITA method is able to filter a surveillance population to generate a subpopulation with a higher density of true positives, regardless of the prevalence and pre-test probability of COVID-19 infection in the population under surveillance for the disease, and, thus, increases testing positivity rates. We also demonstrate the utility of the ITA to filter individuals for allocating diagnostic tests not only in cases of symptomatic individuals but also for asymptomatic individuals who may not be tested and diagnosed otherwise. While the sensitivity and specificity of diagnostic tests are not affected by ITA, this more efficient testing allocation approach identifies more cases in less time and with fewer resources [39]–[42].

The basis of the ITA method is the detection of physiological changes associated with infection onset, which are well-established to be detectable by biometric sensors [18], [22], [24]–[26], [32]–[35], [43], [44]. Consistent with prior literature, we demonstrate here that digital biomarkers derived from heart rate and physical activity are indicative of infection onset. A unique contribution of our work is the demonstration of differences in digital biomarker significance with respect to time prior to the diagnostic test date; specifically, we show that differences in RHR features were significant between COVID-19 positive and negative groups as early as ten days prior to the diagnostic test date whereas differences in most step features were not significant until five days prior to the diagnostic test date. One step feature, minimum  $\Delta$ Steps, was significant up to nine days prior to the diagnostic test date, potentially demonstrating a link between activity levels (and perhaps noncompliance with lockdown measures) and COVID-19 exposure. Furthermore, RHR begins to deviate from baseline earlier than steps (as early as thirteen days vs. ten days prior to the diagnostic test date, respectively), and the peak effect (maximum deviation from baseline) of

infection also occurs earlier in RHR than steps (one day prior vs. two days after the diagnostic test date, respectively) for those who were COVID-19 positive. These results indicate that changes in physiology (RHR) occur earlier in the infection period, while symptoms and reduced physical activity (steps) transpire later in the infection period, when people may limit their movement due either to illness or mandatory quarantine. A recent COVID-19 study assessing prolonged physiological and behavioral changes using wearables also observed that COVID-19 positive individuals took more time to return to their RHR baseline values compared to their step and sleep baseline values following the acute COVID-19 infection period [29]; however, this work explored the post-infection period of the data whereas here we explore the pre-infection period as well as the acute infection period using a systematic grid search approach. Another recent study [32] that developed machine learning models to passively detect COVID-19 using wearable data noted relative changes in feature importance when including data post-diagnosis. However, to our knowledge, we are the first to demonstrate and establish the dynamics of feature importance over time prior to the diagnostic test date, indicating which features should be weighted more heavily in prediction models and when.

Another important contribution of our work is demonstrating the utility of RHR and steps features in the tradeoff between the true positive rate and false positive rate (ROC analysis) and the tradeoff between the true positive rate and the positive predictive value (PR analysis). Specifically, we show that while steps features provide more salient information on the trade-off between the true-positive rate and false positive rate, RHR features provide more salient information on the trade-off between the true positive rate and the precision (positive predictive value). To our knowledge, this is the first demonstration of this tradeoff in predictive model development for COVID-19 infection detection. The ITA model, in addition to using features of RHR and steps, can likely be further extended and improved with features from other digital biomarkers such as skin temperature, respiratory rate, blood oxygen saturation, and sleep duration [21], [22], [33], [34]. It is anticipated that each of these distinct digital biomarkers would capture a physiological response to infection at different times during the detection period, thus improving the robustness and overall performance of the ITA approach.

One of the important observations from our work was the clear separation of the 95% confidence intervals of the means of digital biomarkers between COVID-19 positive and negative populations as early as five days prior to the test date (Fig. 2A and Extended Data Fig. 3A and 4A), while the variances of the groups have overlapping distributions in the same time window (Extended Data Fig. 5). Notably, a lack of overlap in 95% confidence intervals does not necessarily imply significant differences between the groups [30] as standard deviation is a valuable descriptive measure of the data that should be considered as well. There are many possible sources of variance in studies involving wearable data, including the inclusion of different device types and technologies, contexts of measurement (e.g., time of day, activity type, etc.), differences in physiological response to infection, etc. We mitigated this issue by segmenting by device type and data resolution, as well as by utilizing measurements during resting periods only for the RHR calculation. In the future, larger datasets can enable segmentation by demographics (e.g., age, sex, weight, etc.) that would likely further reduce the variance. Sharing datasets between studies, as demonstrated here, can also augment the study population and further reduce the variance. An open

question is whether the resolution of current photoplethysmography-based wearable heart rate technologies is high enough to adequately detect signals above the population variance.

Here, we did not deploy the ITA method in real-time and, thus, its performance in practice still remains to be tested. Both the CovIdentify and MyPHD studies were primarily Bring Your Own Device (BYOD) study designs, in which people who already own smart devices are recruited to participate. The BYOD design presents two major challenges: 1) participants must own a smart device, which limits eligibility to those who can afford devices, and 2) many different types of devices are used, introducing an additional source of noise in the analysis. We mitigated the first challenge by developing and implementing the Demographic Improvement Guideline [45], and the second challenge by dividing our overall dataset into cohorts with homogeneous sampling frequencies and/or device types. Although we recognize that certain factors decrease the likelihood of wearable device ownership, such as lower income or living in a rural area [46]–[48], the precipitously decreasing cost of wearable technology is rapidly increasing the equitable distribution of these technologies [49].

Another limitation of the study is the data missingness and its impact on the deviation of the digital biomarkers, as the source of missingness may confound the disease-related physiological variation. For example, we observed that some participants in our study did not wear their devices when they were feeling sick, as observed in other studies [25], which resulted in a reduction in recorded physical activity. For that reason, it can be a challenge to isolate the effects of physiological and behavioral changes on the digital biomarkers. Furthermore, some devices require more frequent charging (e.g., Apple Watch), which results in more missing data that may also impact model performance. We mitigated this challenge by further developing our model on a single device and homogeneous sampling frequency (FHF) cohort.

The recent body of work on COVID-19 detection using smartwatches uses AUC-ROC to evaluate model performance [26], [32]–[35], which is only an appropriate metric for class-balanced data, and is otherwise misleading [36], [37]. In these large-scale studies conducted on a convenience sample of the population for a disease with low prevalence, there exists an inherent challenge of class imbalance because most of the study population does not contract the disease. This was a challenge that we faced in our study, and, further complicating matters, many of the COVID-19 positive participants did not wear their wearable devices at the start of their infection, exacerbating the class imbalance. While less frequently reported than AUC-ROC, the AUC-PR is the correct evaluation metric for evaluating a classifier on imbalanced data [36], which is what we report here. We show that even with a strong AUC-ROC, the AUC-PR demonstrates the limitations of performance. Methods to resolve class imbalance, especially when working with wearable device data, can be further investigated for future studies.

While our study focused on improving testing allocation for COVID-19, the methods developed herein are extensible to other types of infections and could be used to fortify our future pandemic preparedness. Using ITA to improve disease surveillance could be especially important in underserved communities which may benefit from the fact that the ITA method is useful even with only steps digital biomarkers

which may be obtained from smartphones which are owned by 85% of the population in the US [50] and up to 76% globally [51]. By targeting diagnostic testing toward individuals who are more likely to truly be infected with a disease, we can improve the allocation and utility of diagnostic tests, ultimately reducing mortality and increasing our ability to control the current and future pandemics.

## Methods

### Participant Recruitment and Data Collection

The CovIdentify study launched on April 2, 2020 (Duke University Institutional Review Board #2020 – 0412). Eligibility criteria included age over 18 years and internet access. Social networks and social media advertising were used to recruit participants. By May 25, 2021, a total of 7,348 participants were recruited and e-consented through the Research Electronic Data Capture (REDCap) system [23]. During enrollment, participants were given the option to donate 12 months of retrospective wearable data and 12 months of prospective wearable data. Wearable data was collected via the CovIdentify iOS app for devices connected to the Apple Health kit (e.g., Apple Watch) or via Application Programming Interfaces (APIs) for other devices (e.g., Garmin and Fitbit devices). The participants were also asked to complete an onboarding (enrollment) survey and daily surveys. The surveys were in English or Spanish and included questions on symptoms, social distancing, diagnostic testing results, and related information (Supplementary Document 1). Surveys were collected using the CovIdentify iOS app, text messaging, and/or emails. All wearable data and survey results were stored in a secured Microsoft Azure data platform and later analyzed in the Microsoft Azure Machine Learning environment. Soon after CovIdentify was launched, exploratory data analysis revealed major differences between CovIdentify demographics and the demographics of COVID-19 positive cases and deaths in the U.S., as well as overall U.S. demographics based on the 2020 U.S. Census [52], [53]. We sought to mitigate the imbalance throughout the duration of the study by providing wearable devices to underrepresented populations [45]. COVID-19 vaccine reporting was added to the daily surveys in February 2021, where we asked questions regarding the vaccination date, vaccine brand, vaccine related symptoms, and dose number.

Wearable data processing and analysis: Participants were asked to fill out an enrollment survey following the informed e-consent. Daily symptom survey and wearable data from the participants were analyzed both separately and together. For the overall analysis, we only included participants with self-reported diagnostic test results for COVID-19. These participants were further divided into two categories based on the self-reported diagnostic test results: COVID-19 positive and COVID-19 negative.

In addition to the data collected via CovIdentify, we augmented our analysis by including data from the MyPHD study, as reported on in the two recent publications by Mishra et al. [25] and Alavi et al. [24]. The data from Mishra et al. included heart rate, step count, and sleep data for 27 COVID-19 positive cases. It also included metadata of symptom onset and test dates. The data from Alavi et al. included heart rate and step count data for 83 COVID-19 positive cases and 1,019 COVID-19 negative cases as well as metadata including symptom onset and test dates.

For wearable data analysis, we only included days of wearable data when both heart rate and step count were available. Out of the 1,239 participants (113 from CovIdentify and 1,126 from MyPhD study) who had both heart rate and step count data available, we had device-reported daily values of RHR and step count for 67 participants, and high frequency (second- or minute- level, depending on device types) wearable data for 1,172 participants. For participants with high frequency heart rate data, we calculated daily RHR from the heart rate data points recorded between midnight and 7 AM, when there were no steps recorded. For those participants with available high frequency wearable data, we chose a data-driven threshold (i.e., a minimum number of heart rate data points between midnight and 7 AM with zero recorded steps) to include our calculated RHR data from that day in the subsequent analysis. As the sampling rate varies by device types (Fitbit, Garmin, and Apple Watch), we generated separate data distributions of the datasets for these three device types and selected the first quartile of heart rate data points per device as the data-driven threshold, which resulted in a threshold of 2,630, 19, and 1,389 heart rate data points for Fitbit, Apple Watch, and Garmin devices, respectively. In other words, on a given day, a participant with Fitbit wearable data needed to have at least 2,630 heart rate data points between midnight and 7 AM with zero recorded steps for us to include our calculated RHR value in the subsequent analysis. Following this intraday data point threshold, we used an interday data threshold: a minimum number of days with available wearable data to be included in the analysis (50% in the baseline period and 50% between nine days and one day prior to the diagnostic test date in the detection period). We explored different minimum number of days of available wearable data in the baseline and detection periods and selected these two thresholds to maximize the number of participants while keeping the performance of the ITA model on the training dataset consistent, defined as less than 10% variation of the performance metrics (AUC-ROC and AUC-PR)).

Cohort Definition: The wearable data availability thresholds (both intraday and interday) resulted in an AF cohort of 520 participants (83 from CovIdentify and 437 from MyPHD) with sufficient wearable data. We then created two more subsets from this cohort (Extended Data Fig. 2): (1) AHF cohort: participants with high frequency wearable data (469 participants, 54 COVID-19 positive and 415 COVID-19 negative), and (2) FHF cohort: participants with high frequency wearable data from a single source (Fitbit) (280 participants, 40 COVID-19 positive and 240 COVID-19 negative) to explore the impact of utilizing wearable data from different sources and resolutions on the ITA model development. We employed these three cohorts separately for the ITA model development and compared the resulting models' performance in the corresponding training and test datasets of these cohorts. We divided each cohort into an 80% train and 20% test split, with FHF as a subgroup of AHF (which itself is a subset of AF) to ensure that no observations in the training dataset of one cohort existed in the test dataset of another (Extended Data Fig. 2).

## Digital biomarker definition:

Given the use of datasets with different device types, a consistent RHR definition was used in order to harmonize the cohorts with high frequency wearable data. We calculated the daily RHR digital biomarker by aggregating the high frequency heart rate data points available between midnight and 7 AM, when

there were no steps recorded. Step count was calculated by summing all recorded step values during a 24 hour period in order to produce a daily step count digital biomarker.

## **Feature Engineering and Extraction**

Following the creation of three cohorts (AF, AHF, and FHF) and their corresponding training and test sets, we performed exploratory data analysis (EDA) and extracted features from the time-series digital biomarkers (RHR and step count). For the EDA on the time-series digital biomarkers, we explored the difference in trajectories of digital biomarkers between COVID-19 positive and COVID-19 negative participants (Fig. 2A and Extended Data Fig. 3A and 4A). Following the EDA, we extracted the features mentioned in Table 2 from the raw digital biomarkers. We first calculated four deviation metrics, which capture the deviation in digital biomarkers from participants' baseline during the detection phase. Following the deviation metrics calculation, we calculated summary statistics of these four deviation metrics which we refer as to features for this manuscript. We extracted the same features from the training and test datasets. Following the feature extraction, we performed statistical analysis on the features from the training datasets of the three cohorts to see which features are statistically different between the two groups and how their significance levels vary with different detection period combinations (detection end date (DED) and duration of the detection window (DWL)) using a systematic grid search to optimize DED and DWL (Fig. 2B and Extended Data Fig. 3B and 4B). We utilized multiple hypothesis testing with Benjamini-Hochberg adjusted p-values for this statistical analysis. Following the statistical analysis and systematic grid search to obtain the optimal detection period to extract the features, we only utilized the intersection of the statistically significant features ( $p\text{-value} < 0.05$ ) extracted from digital biomarkers recorded between five days and one day and three days and one day prior to the diagnostic test date for the development of the ITA model.

## **ITA Model Development**

Following feature extraction, we developed predictive models to classify COVID-19 positive and negative participants in the training dataset of each cohort (AF, AHF, and FHF) using nested cross-validation (CV) and later validated the models on corresponding independent test datasets. We chose five state-of-the-art machine learning models (logistic regression, K-nearest neighbor (KNN), support vector machine (SVM), Random Forest (RF), and extreme gradient boosting (XGBoost) [54], [55]) for the development of the ITA models to explore how increasing model complexity and adding non-linearity would impact the model performance. We trained these classification models on the training dataset using nested CV with an inner CV loop for hyperparameter tuning and an outer CV loop for model selection. For model training, we selected recall as our preferred scoring metric for model selection to emphasize the relative impact/cost of false negatives compared to false positives, as an individual who is truly positive for COVID-19 and is wrongly classified as negative (or healthy) would further spread disease. For model performance evaluation, we used two well-established reporting metrics, including the most commonly reported metric for studies of this kind (the area under the curve for the receiver operating characteristic curve (AUC-ROC)) [26], [31]–[35], and the metric that is most appropriate for this classification task (AUC for the precision-

recall curve (AUC-PR)) [36] (Extended Data Table 1; Fig. 3, Fig. 4, and Extended Data Fig. 6). AUC-PR is more appropriate with class-imbalanced data [36], [37], which is the case here (12–15% COVID-19 positive and 85–88% negative for each of the three cohorts). The results reported for the training dataset (Extended Data Table 1, Fig. 3A, B, and C, Fig. 4A, B, and C, and Extended Data Fig. 6A, B, and C) were generated from the validation on the held-out dataset (fold) from each iteration of the outer CV loop which was not used in the model training. Based on the CV results of the five machine learning models on the training dataset, we chose the logistic regression model to further evaluate performance on the independent testing dataset (Fig. 3D, E, and F, Fig. 4D, E, and F, and Extended Data Fig. 6D, E, and F). For validation on the independent test dataset, we trained the logistic regression model on the entire training dataset using a grid search with five stratified folds for hyperparameter tuning and selected the best model (with tuned hyperparameters) to validate on the test dataset.

### **Nested-Cross Validation**

For model development with the training dataset, we utilized nested CV over as traditional CV, which is a common approach in similar studies [26], [32], [34], [35], because it uses the same data for hyperparameter tuning and model performance evaluation [56]. In nested CV (also called double CV), the hyperparameter tuning procedure is nested (inner loop) under the model selection procedure (outer loop) and the inner loop is used for optimizing the hyperparameters of the model with inner CV, and the outer loop is used to compute the error of the optimized model with outer CV. [57]. For the nested CV, we divided the training set into ten stratified folds (keeping the ratio of COVID-19 positive and negative participants the same across each fold) for the outer loop. For each iteration of the outer loop, the model was trained on data from nine folds by optimizing the hyperparameters of the model with inner CV, and validating on the left-out fold, a process which was repeated nine more times. In each iteration of the outer loop, the outer training data (from nine folds) were further divided into five stratified folds (inner loop) to tune hyperparameters using a grid search. Out of the five iterations with the grid search in the inner loop, the best model (including hyperparameters) was selected, and this model was used in the model performance evaluation in the outer loop. This way of model development using two CV steps separates hyperparameter tuning and model selection in order to reduce bias in model performance.

### **Feature Importance Ranking**

To calculate the feature importance ranking, we trained the logistic regression model using a grid search with five stratified folds for hyperparameter tuning and selected the best model (with optimized hyperparameters) to train on the entire training set of each cohort and extracted the coefficients for each feature used in the optimized model. We reported the absolute value of each coefficient as the relative importance of the features (Extended Data Fig. 8).

## **Declarations**

## **Acknowledgements.**

This work was supported by in part by Duke OIT, Duke Bass Connections Fellowship, Duke Margolis Center for Health Policy, Duke MEDx, Microsoft AI in Health, Duke CTSI (*UL1TR002553*), and NC Biotech (2020-FLG-3884). This article was prepared while Geoffrey S Ginsburg was employed at Duke University. Dana K Pasquale was funded by NIH/NICHD (R25HD079352) and CDC (BAA 75D301-20-R-68024). The opinions expressed in this article are the author's own and do not reflect the view of the National Institutes of Health, the Department of Health and Human Services, or the United States government. We would like to also thank Leatrice Martin, Veronica Palacios-Grandez, John Owens, Julie Ekstrand, Cecilia Plez, Hugh Thomas, Philabian Lindo, Richard Outten, Shellene Walker, Tracey Futhey, Jimmy Dorff, Rob Carter, Sean Dilda, Vanessa Simmons, Andy Ingham, Charley Kneifel, Andrew Olson, Whitney Welsh, Jonathan McCall, Margaret Pendzich, Marialuisa Solis-Guzman, Erich Huang, and Victoria Christian for their contributions to the infrastructure and their dedication to this project.

## Tables

**Table 1. Summary of the Cohorts.** Total refers to training + test data

Cohort	Total N (Test N)	Total COVID+ (Test)	Total COVID- (Test)
<b>All-Frequency (AF)</b>	520 (105)	63 (13)	457 (92)
<b>All-High-Frequency (AHF)</b>	469 (97)	54 (11)	415 (86)
<b>Fitbit-High-Frequency (FHF)</b>	280 (63)	40 (7)	240 (56)

**Table 2:** Features Extracted from the Digital Biomarkers (DBs) for the Development of ITA Algorithm

Metric	Definition	Equation
<b>Deviation Metrics</b>		
Delta ( $\Delta$ )	Deviation in digital biomarker from baseline median value	$DB_{\text{Detection}} - DB_{\text{Baseline, Median}}$
Delta_Normalized	Delta normalized by baseline median value	$((DB_{\text{Detection}} - DB_{\text{Baseline, Median}}) / DB_{\text{Baseline, Median}})$
Delta_Standardized	Delta standardized by baseline median and interquartile range (IQR)	$((DB_{\text{Detection}} - DB_{\text{Baseline, Median}}) / DB_{\text{Baseline, IQR}})$
Z-score	Deviation in digital biomarker from baseline mean, standardized by baseline standard deviation (SD)	$((DB_{\text{Detection}} - DB_{\text{Baseline, Mean}}) / DB_{\text{Baseline, SD}})$
<b>Summary Statistics (Features)</b>		
Average	Average of inter-day deviation metrics	
Median	Median of inter-day deviation metrics	
Maximum	Maximum of inter-day deviation metrics	
Minimum	Minimum of inter-day deviation metrics	
Range	Range of inter-day deviation metrics	

## References

- [1] "COVID Live - Coronavirus Statistics - Worldometer." <https://www.worldometers.info/coronavirus/> (accessed Mar. 25, 2022).
- [2] "United States COVID - Coronavirus Statistics - Worldometer." <https://www.worldometers.info/coronavirus/country/us/> (accessed Mar. 25, 2022).
- [3] S. M. Moghadas *et al.*, "The implications of silent transmission for the control of COVID-19 outbreaks," *Proc. Natl. Acad. Sci.*, vol. 117, no. 30, pp. 17513–17515, Jul. 2020, doi: 10.1073/pnas.2008373117.
- [4] "Why Pretest and Posttest Probability Matter in the Time of COVID-19," *ASM.org*. <https://asm.org/Articles/2020/June/Why-Pretest-and-Posttest-Probability-Matter-in-the> (accessed Feb. 22, 2022).
- [5] S. G. Stolberg and S. LaFraniere, "With Omicron, U.S. Testing Capacity Faces Intense Pressure," *The New York Times*, Dec. 18, 2021. Accessed: Feb. 24, 2022. [Online]. Available: <https://www.nytimes.com/2021/12/17/us/politics/us-covid-tests-omicron.html>
- [6] "Reports of price gouging amid shortages of COVID-19 tests." <https://www.cbsnews.com/news/at-home-test-covid-price-gouging/> (accessed Feb. 24, 2022).

- [7] C. O'donnell, "U.S. COVID-19 tests again in short supply as infections soar, schools reopen," *Reuters*, Aug. 27, 2021. Accessed: Feb. 22, 2022. [Online]. Available: <https://www.reuters.com/world/us/us-covid-19-tests-again-short-supply-infections-soar-schools-reopen-2021-08-27/>
- [8] "Omicron testing shortages and delays are making results useless—and deepening COVID inequality," *Fortune*. <https://fortune.com/2022/01/10/omicron-testing-shortages-delays-covid-inequality/> (accessed Feb. 24, 2022).
- [9] R. Heilweil, "How omicron broke Covid-19 testing," *Vox*, Dec. 21, 2021. <https://www.vox.com/recode/2021/12/21/22848286/omicron-rapid-test-covid-19-antigen> (accessed Feb. 24, 2022).
- [10] "More coronavirus tests will be available next month, Fauci says, as U.S. struggles with shortage," *Washington Post*. Accessed: Feb. 24, 2022. [Online]. Available: <https://www.washingtonpost.com/health/2021/12/27/omicron-covid-test-shortage-fauci/>
- [11] B. Rader *et al.*, "Geographic access to United States SARS-CoV-2 testing sites highlights healthcare disparities and may bias transmission estimates," *J. Travel Med.*, vol. 27, no. 7, p. taaa076, Nov. 2020, doi: 10.1093/jtm/taaa076.
- [12] J. M. Souch and J. S. Cossman, "A Commentary on Rural-Urban Disparities in COVID-19 Testing Rates per 100,000 and Risk Factors," *J. Rural Health*, vol. 37, no. 1, pp. 188–190, 2021, doi: 10.1111/jrh.12450.
- [13] "New Surgo Analysis Identifies Highly Vulnerable Rural Communities as COVID-19 Testing Deserts," *Surgo Ventures*. <https://surgoventures.org/portfolio/action-areas/new-surgo-analysis-identifies-highly-vulnerable-rural-communities-as-covid-19-testing-deserts> (accessed Feb. 24, 2022).
- [14] S. Kanno, S. Kandula, and J. Shaman, "The association between early country-level COVID-19 testing capacity and later COVID-19 mortality outcomes," *Influenza Other Respir. Viruses*, vol. 16, no. 1, pp. 56–62, 2022, doi: 10.1111/irv.12906.
- [15] CDC, "Community, Work, and School," *Centers for Disease Control and Prevention*, Feb. 11, 2020. <https://www.cdc.gov/coronavirus/2019-ncov/community/health-equity/racial-ethnic-disparities/increased-risk-exposure.html> (accessed Jun. 13, 2021).
- [16] J. M. Radin, N. E. Wineinger, E. J. Topol, and S. R. Steinhubl, "Harnessing wearable device data to improve state-level real-time surveillance of influenza-like illness in the USA: a population-based study," *Lancet Digit. Health*, vol. 2, no. 2, pp. e85–e93, 2020.
- [17] X. Li *et al.*, "Digital health: tracking physiomes and activity using wearable biosensors reveals useful health-related information," *PLoS Biol.*, vol. 15, no. 1, p. e2001402, 2017.

- [18] A. Shapiro *et al.*, "Characterizing COVID-19 and Influenza Illnesses in the Real World via Person-Generated Health Data," *Patterns*, vol. 2, no. 1, p. 100188, 2021.
- [19] T. Mishra *et al.*, "Pre-symptomatic detection of COVID-19 from smartwatch data," *Nat. Biomed. Eng.*, vol. 4, no. 12, pp. 1208–1220, 2020.
- [20] G. Quer *et al.*, "Wearable sensor data and self-reported symptoms for COVID-19 detection," *Nat. Med.*, vol. 27, no. 1, pp. 73–77, 2021.
- [21] D. J. Miller *et al.*, "Analyzing changes in respiratory rate to predict the risk of COVID-19 infection," *PLoS One*, vol. 15, no. 12, p. e0243693, 2020.
- [22] E. Grzesiak *et al.*, "Assessment of the Feasibility of Using Noninvasive Wearable Biometric Monitoring Sensors to Detect Influenza and the Common Cold Before Symptom Onset," *JAMA Netw. Open*, vol. 4, no. 9, p. e2128534, Sep. 2021, doi: 10.1001/jamanetworkopen.2021.28534.
- [23] P. A. Harris, R. Taylor, R. Thielke, J. Payne, N. Gonzalez, and J. G. Conde, "Research Electronic Data Capture (REDCap) - A metadata-driven methodology and workflow process for providing translational research informatics support," *J. Biomed. Inform.*, vol. 42, no. 2, pp. 377–381, Apr. 2009, doi: 10.1016/j.jbi.2008.08.010.
- [24] A. Alavi *et al.*, "Real-time alerting system for COVID-19 and other stress events using wearable data," *Nat. Med.*, pp. 1–10, Nov. 2021, doi: 10.1038/s41591-021-01593-2.
- [25] T. Mishra *et al.*, "Pre-symptomatic detection of COVID-19 from smartwatch data," *Nat. Biomed. Eng.*, pp. 1–13, Nov. 2020, doi: 10.1038/s41551-020-00640-6.
- [26] G. Quer *et al.*, "Wearable sensor data and self-reported symptoms for COVID-19 detection," *Nat. Med.*, pp. 1–5, Oct. 2020, doi: 10.1038/s41591-020-1123-x.
- [27] G. Rapeport, E. Smith, A. Gilbert, A. Catchpole, H. McShane, and C. Chiu, "SARS-CoV-2 Human Challenge Studies – Establishing the Model during an Evolving Pandemic," *N. Engl. J. Med.*, vol. 385, no. 11, pp. 961–964, Sep. 2021, doi: 10.1056/NEJMp2106970.
- [28] CDC, "COVID-19 and Your Health," *Centers for Disease Control and Prevention*, Feb. 11, 2020. <https://www.cdc.gov/coronavirus/2019-ncov/symptoms-testing/testing.html> (accessed Mar. 09, 2022).
- [29] J. M. Radin *et al.*, "Assessment of Prolonged Physiological and Behavioral Changes Associated With COVID-19 Infection," *JAMA Netw. Open*, vol. 4, no. 7, pp. e2115959–e2115959, Jul. 2021, doi: 10.1001/jamanetworkopen.2021.15959.
- [30] M. Krzywinski and N. Altman, "Error bars," *Nat. Methods*, vol. 10, no. 10, Art. no. 10, Oct. 2013, doi: 10.1038/nmeth.2659.

- [31] A. P. Bradley, "The use of the area under the ROC curve in the evaluation of machine learning algorithms," *Pattern Recognit.*, vol. 30, no. 7, pp. 1145–1159, Jul. 1997, doi: 10.1016/S0031-3203(96)00142-2.
- [32] M. Gadaleta *et al.*, "Passive detection of COVID-19 with wearable sensors and explainable machine learning algorithms," *Npj Digit. Med.*, vol. 4, no. 1, pp. 1–10, Dec. 2021, doi: 10.1038/s41746-021-00533-1.
- [33] A. E. Mason *et al.*, "Detection of COVID-19 using multimodal data from a wearable device: results from the first TemPredict Study," *Sci. Rep.*, vol. 12, no. 1, Art. no. 1, Mar. 2022, doi: 10.1038/s41598-022-07314-0.
- [34] B. Conroy *et al.*, "Real-time infection prediction with wearable physiological monitoring and AI to aid military workforce readiness during COVID-19," *Sci. Rep.*, vol. 12, no. 1, Art. no. 1, Mar. 2022, doi: 10.1038/s41598-022-07764-6.
- [35] A. Natarajan, H.-W. Su, and C. Heneghan, "Assessment of physiological signs associated with COVID-19 measured using wearable devices," *Npj Digit. Med.*, vol. 3, no. 1, pp. 1–8, Nov. 2020, doi: 10.1038/s41746-020-00363-7.
- [36] T. Saito and M. Rehmsmeier, "The Precision-Recall Plot Is More Informative than the ROC Plot When Evaluating Binary Classifiers on Imbalanced Datasets," *PLoS ONE*, vol. 10, no. 3, p. e0118432, Mar. 2015, doi: 10.1371/journal.pone.0118432.
- [37] J. Davis and M. Goadrich, "The relationship between Precision-Recall and ROC curves," in *Proceedings of the 23rd international conference on Machine learning - ICML '06*, Pittsburgh, Pennsylvania, 2006, pp. 233–240. doi: 10.1145/1143844.1143874.
- [38] J. Lever, M. Krzywinski, and N. Altman, "Model selection and overfitting," *Nat. Methods*, vol. 13, no. 9, Art. no. 9, Sep. 2016, doi: 10.1038/nmeth.3968.
- [39] I. of M. (US) C. on H. C. Technology, H. Sox, S. Stern, D. Owens, and H. L. Abrams, *The Use of Diagnostic Tests: A Probabilistic Approach*. National Academies Press (US), 1989. Accessed: Feb. 24, 2022. [Online]. Available: <https://www.ncbi.nlm.nih.gov/books/NBK235178/>
- [40] "Frequently Asked Questions (FAQs) | grants.nih.gov." <https://grants.nih.gov/faqs#/inclusion-basic-on-sex-gender-and-race-ethnicity.htm> (accessed Jul. 06, 2021).
- [41] J. Watson, P. F. Whiting, and J. E. Brush, "Interpreting a covid-19 test result," *BMJ*, vol. 369, p. m1808, May 2020, doi: 10.1136/bmj.m1808.
- [42] J. B. Gubbay *et al.*, "Impact of coronavirus disease 2019 (COVID-19) pre-test probability on positive predictive value of high cycle threshold severe acute respiratory coronavirus virus 2 (SARS-CoV-2) real-time reverse transcription polymerase chain reaction (RT-PCR) test results," *Infect. Control Hosp. Epidemiol.*, pp. 1–5, doi: 10.1017/ice.2021.369.

- [43] X. Li *et al.*, “Digital Health: Tracking Physiomes and Activity Using Wearable Biosensors Reveals Useful Health-Related Information,” *PLOS Biol.*, vol. 15, no. 1, p. e2001402, Jan. 2017, doi: 10.1371/journal.pbio.2001402.
- [44] J. M. Radin, N. E. Wineinger, E. J. Topol, and S. R. Steinhubl, “Harnessing wearable device data to improve state-level real-time surveillance of influenza-like illness in the USA: a population-based study,” *Lancet Digit. Health*, vol. 2, no. 2, pp. e85–e93, Feb. 2020, doi: 10.1016/S2589-7500(19)30222-5.
- [45] P. J. Cho *et al.*, “The Demographic Improvement Guideline to Reduce Bias Resulting from Bring-Your-Own-Device Study Design (Preprint),” *JMIR mHealth and uHealth*, preprint, Apr. 2021. doi: 10.2196/preprints.29510.
- [46] E. a Vogels, “Digital divide persists even as Americans with lower incomes make gains in tech adoption,” *Pew Research Center*. <https://www.pewresearch.org/fact-tank/2021/06/22/digital-divide-persists-even-as-americans-with-lower-incomes-make-gains-in-tech-adoption/> (accessed Mar. 21, 2022).
- [47] E. a Vogels, “Some digital divides persist between rural, urban and suburban America,” *Pew Research Center*. <https://www.pewresearch.org/fact-tank/2021/08/19/some-digital-divides-persist-between-rural-urban-and-suburban-america/> (accessed Mar. 21, 2022).
- [48] S. Macridis, N. Johnston, S. Johnson, and J. K. Vallance, “Consumer physical activity tracking device ownership and use among a population-based sample of adults,” *PLoS ONE*, vol. 13, no. 1, p. e0189298, Jan. 2018, doi: 10.1371/journal.pone.0189298.
- [49] K. Guk *et al.*, “Evolution of Wearable Devices with Real-Time Disease Monitoring for Personalized Healthcare,” *Nanomaterials*, vol. 9, no. 6, p. 813, May 2019, doi: 10.3390/nano9060813.
- [50] 1615 L. St NW, Suite 800 Washington, and D. 20036 USA 202-419-4300 | M.-857-8562 | F.-419-4372 | M. Inquiries, “Demographics of Mobile Device Ownership and Adoption in the United States,” *Pew Research Center: Internet, Science & Tech*. <https://www.pewresearch.org/internet/fact-sheet/mobile/> (accessed Jun. 14, 2020).
- [51] L. Silver, “Smartphone Ownership Is Growing Rapidly Around the World, but Not Always Equally,” *Pew Research Center’s Global Attitudes Project*, Feb. 05, 2019. <https://www.pewresearch.org/global/2019/02/05/smartphone-ownership-is-growing-rapidly-around-the-world-but-not-always-equally/> (accessed Mar. 23, 2022).
- [52] “U.S. Census Bureau QuickFacts: United States.” <https://www.census.gov/quickfacts/fact/table/US/RHI125219> (accessed Jul. 15, 2020).
- [53] CDC, “COVIDView, Key Updates for Week 23,” *Centers for Disease Control and Prevention*, Jun. 12, 2020. <https://www.cdc.gov/coronavirus/2019-ncov/covid-data/covidview/index.html> (accessed Jun. 16, 2020).

[54] S. Uddin, A. Khan, M. E. Hossain, and M. A. Moni, "Comparing different supervised machine learning algorithms for disease prediction," *BMC Med. Inform. Decis. Mak.*, vol. 19, p. 281, Dec. 2019, doi: 10.1186/s12911-019-1004-8.

[55] R. C. Deo, "Machine Learning in Medicine," *Circulation*, vol. 132, no. 20, pp. 1920–1930, Nov. 2015, doi: 10.1161/CIRCULATIONAHA.115.001593.

[56] G. C. Cawley and N. L. C. Talbot, "On Over-fitting in Model Selection and Subsequent Selection Bias in Performance Evaluation," p. 29.

[57] S. Varma and R. Simon, "Bias in error estimation when using cross-validation for model selection," *BMC Bioinformatics*, vol. 7, p. 91, Feb. 2006, doi: 10.1186/1471-2105-7-91.

## Figures

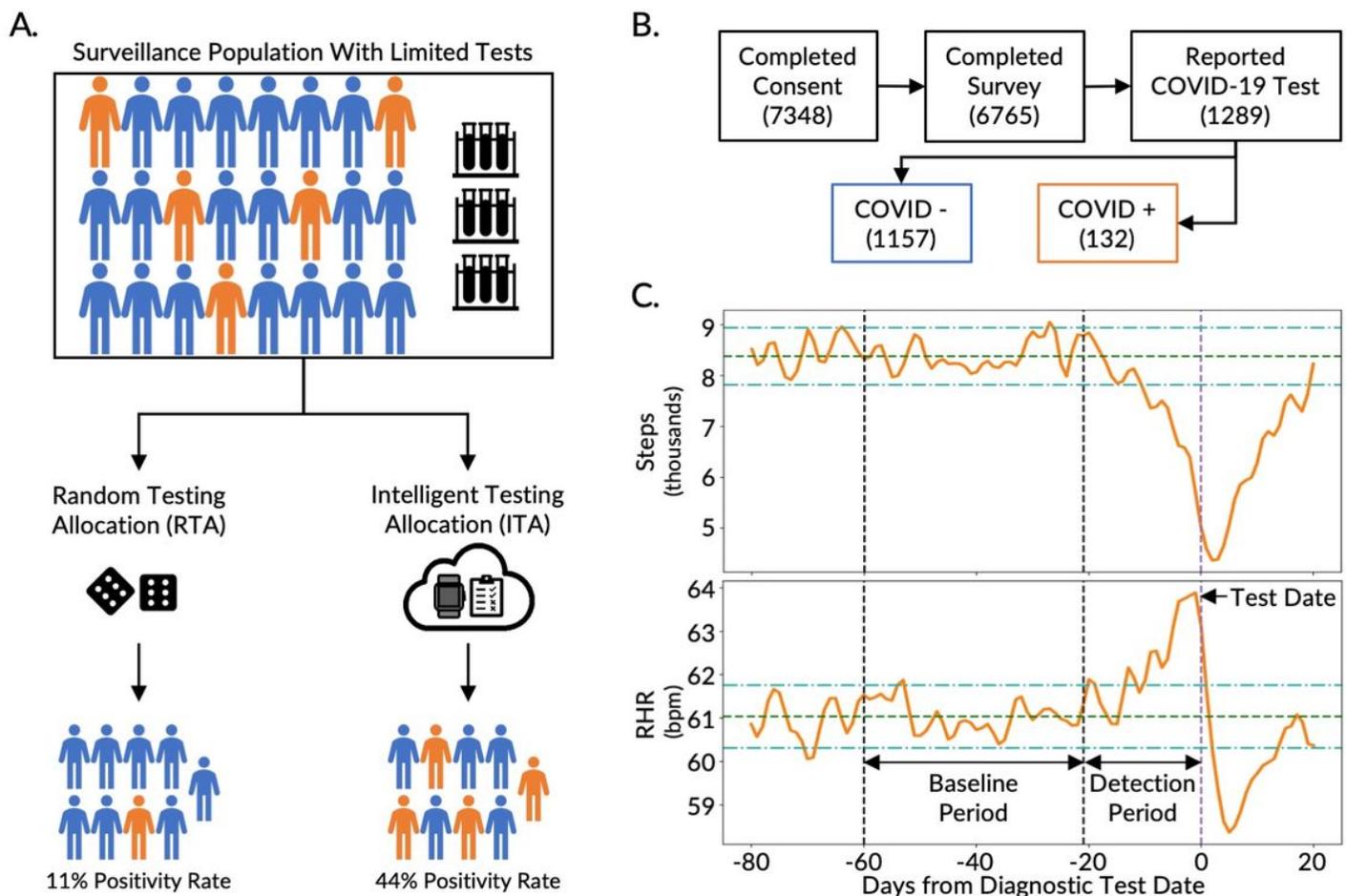
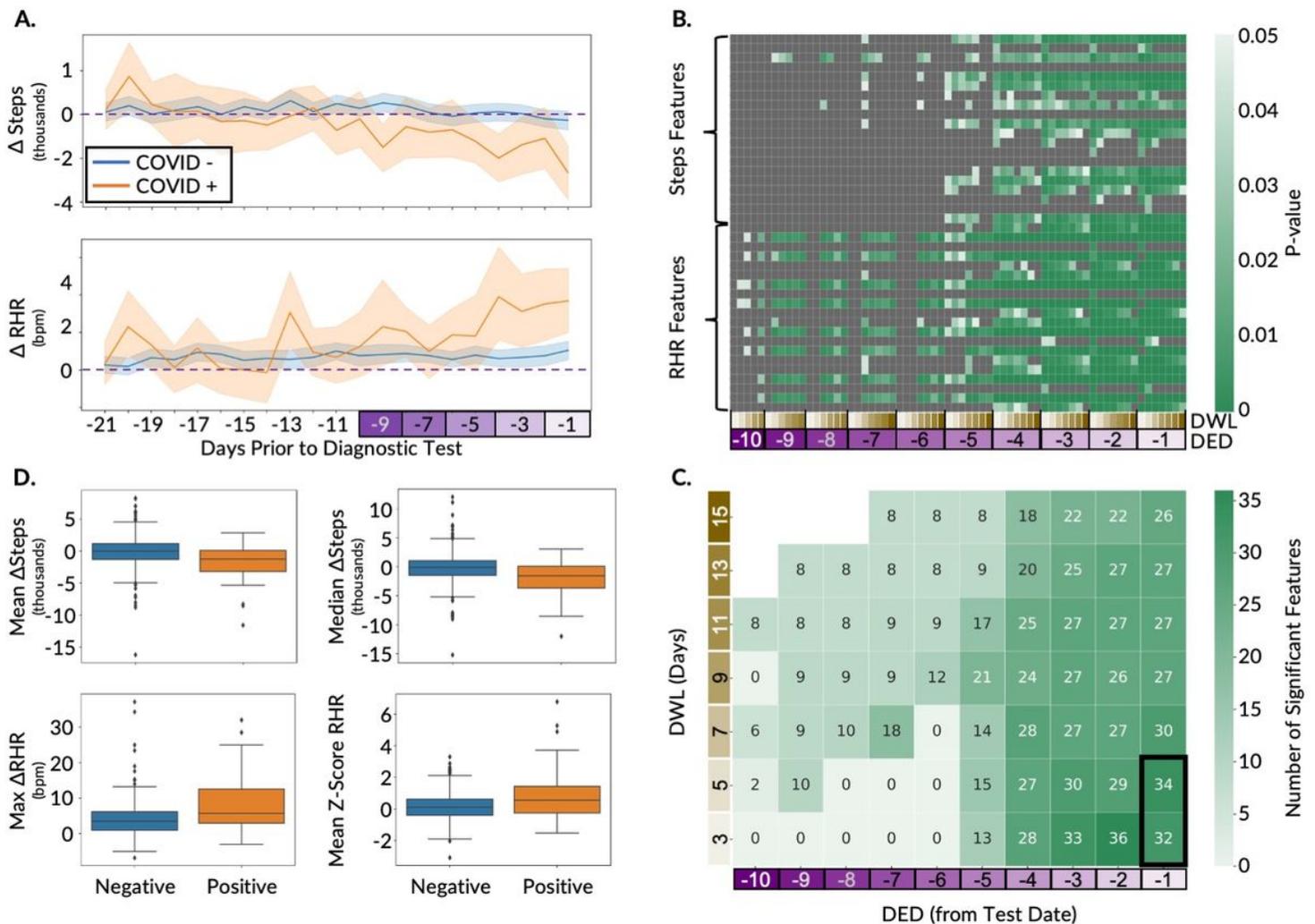


Figure 1

**Overview of the Intelligent Testing Allocation (ITA) model, and the CovIdentify cohort, and data. A.** Overview of the ITA model in comparison to a Random Testing Allocation (RTA) model which demonstrates the benefit of using the ITA model over existing RTA methods to improve the positivity rate

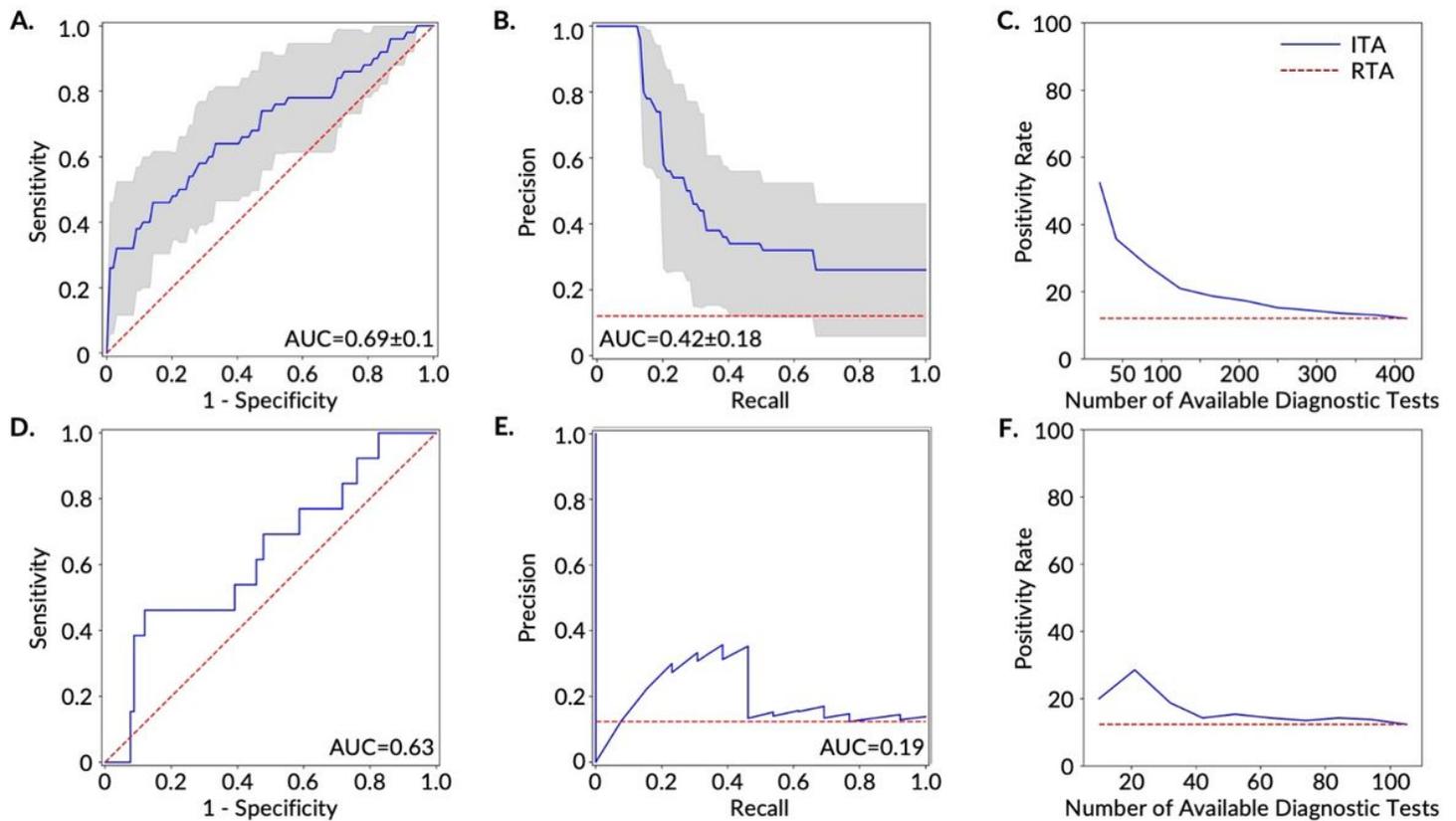
of diagnostic testing in resource-limited settings. **B.** A total of 7,348 participants were recruited following informed consent in the CovIdentify study, out of which 1,289 participants reported COVID-19 diagnostic tests (1,157 diagnosed as negative for COVID-19 and 132 diagnosed as positive for COVID-19). **C.** The top panel shows the time-averaged step count and the bottom panel shows the time-averaged resting heart rate (RHR) of all participants (n=50) in the training set (Extended Data Fig 2, blue) who tested positive for COVID-19 with the pre-defined baseline (between -60 and -22 days from the diagnostic test) and detection (between -21 and -1 days from the diagnostic test) periods marked with vertical black dashed lines. The dark green dashed lines and the light green dash-dotted lines display the baseline period mean and  $\pm 2$  standard deviations from the baseline mean respectively. The light purple dashed vertical line shows the diagnostic test date.



**Figure 2**

**Overview of digital biomarker exploration and feature engineering for the ITA model development on the AF cohort.** **A.** Time-series plot of the deviation in digital biomarkers ( $\Delta$ Steps and  $\Delta$ RHR) in the detection window compared to baseline periods, between the participants diagnosed as COVID-19 positive and negative. The horizontal dashed line displays the baseline median and the confidence bounds show the 95% confidence intervals. **B.** Heatmaps of steps and RHR features that are statistically significantly

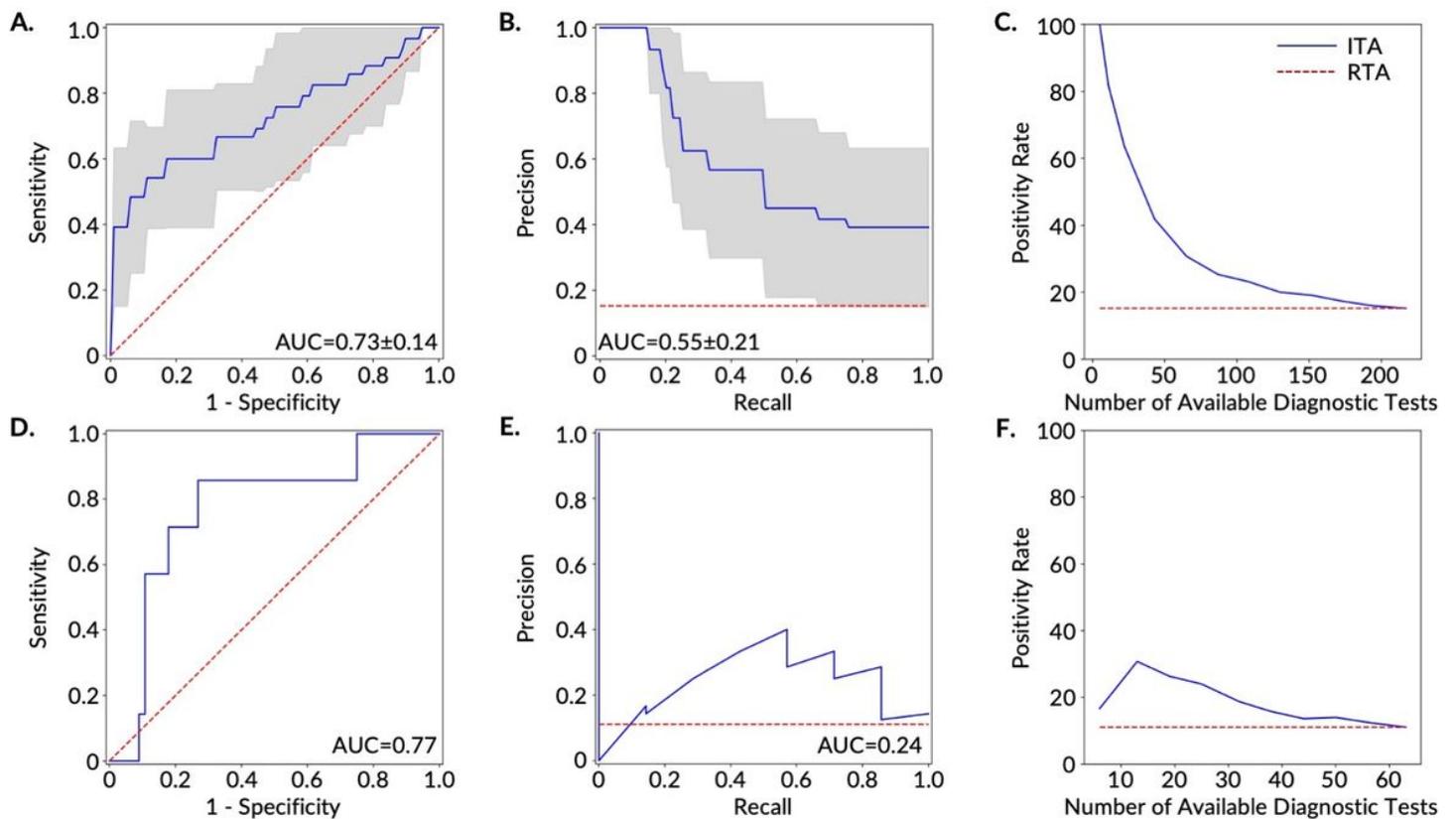
different ( $p$ -value  $< 0.05$ ) in a grid search with different DED and DWL combinations, with green boxes showing  $p$ -values  $< 0.05$  and gray boxes showing  $p$ -values  $\geq 0.05$ . The  $p$ -values are adjusted with the Benjamini-Hochberg method for multiple hypothesis correction. **C.** Summary of the significant features ( $p$ -value  $< 0.05$ ) from B, with each box showing the number of statistically significant features for the different combinations of DED and DWL. The intersection of the significant features across DWL of 3 and 5 days with a common DED of 1 day prior to the test date (as shown using the black rectangle) were used for the ITA model development. **D.** Box plots comparing the distribution of the two most significant steps and RHR features between the participants diagnosed as COVID-19 positive and negative.



**Figure 3**

**Prediction and ranking results of the ITA models on both the training (A-C) and test sets (D-F) for the AF cohort.** **A.** Receiver operating characteristics curves (ROCs) and **B.** precision recall curves (PRCs) for the discrimination between COVID-19 positive participants ( $n=50$ ) and negative participants ( $n=365$ ) in the training set. The gray area shows one standard deviation from the mean of the ROCs/PRCs generated from 10-fold nested cross-validation on the training set and the red dashed line shows the results based on a Random Testing Allocation (RTA) model (the null model). **C.** The positivity rate of the diagnostic testing subpopulation as determined by ITA given a specific number of available diagnostic tests. The red dashed line displays the positivity rate/pre-test probability of an RTA (null) model. **D.** ROC and **E.** PRC for the discrimination between COVID-19 positive participants ( $n=13$ ) and negative participants ( $n=92$ ) in the test set. The red dashed line shows the results based on an RTA model. **F.** Positivity rate of the diagnostic

testing subpopulation as determined by ITA given a specific number of available diagnostic tests. The red dashed line shows the positivity rate of an RTA (null) model.



**Figure 4**

**Prediction and ranking results of the ITA models on both the training (A-C) and test sets (D-F) for the participants with FHF wearable data.** **A.** Receiver operating characteristics curves (ROCs) and **B.** precision-recall curves (PRCs) for the discrimination between COVID-19 positive participants (n=33) and negative participants (n=184) in the training set. The gray area shows one standard deviation from the mean of the ROCs/PRCs generated from 10-fold nested cross-validation on the training set and the red dashed line shows the results based on a Random Testing Allocation (RTA) model (the null model). **C.** The positivity rate of the diagnostic testing subpopulation as determined by ITA given a specific number of available diagnostic tests. The red dashed line displays the positivity rate/pre-test probability of an RTA (null) model. **D.** ROC and **E.** PRC for the discrimination between Covid-19 positive participants (n=7) and negative participants (n=56) in the test set. The red dashed line shows the results based on an RTA model. **F.** Positivity rate of the diagnostic testing subpopulation as determined by ITA given a specific number of available diagnostic tests. The red dashed line shows the positivity rate of an RTA (null) model.

## Supplementary Files

This is a list of supplementary files associated with this preprint. Click to download.

- [SupplementaryDocument1CovIdentifySurvey.pdf](#)
- [SupplementaryMaterialsFiguresandTables.docx](#)
- [ExtendedData.docx](#)