1 Nowcasting epidemic trends using hospital- and community-based vi-

2 rologic test data

- 3 Tse Yang Lim^{1*}, Sanjat Kanjilal^{2,3}, Shira Doron⁴, Jessica Penney⁴, Meredith Haddix⁵, Tae Hee
- 4 Koo⁵, Phoebe Danza⁵, Rebecca Fisher⁵, Yonatan H. Grad^{1,6†}, James A. Hay^{1,7*†}
- 5 *Correspondence to: tseyanglim@hsph.harvard.edu, james.hay@ndm.ox.ac.uk
- ⁶ [†]These authors jointly supervised the work.
- ⁷ ¹ Center for Communicable Disease Dynamics, Harvard T.H. Chan School of Public Health,
- 8 Boston, MA
- ⁹ ² Department of Population Medicine, Harvard Pilgrim Health Care Institute, Boston, MA
- 10 ³ Department of Infectious Diseases, Brigham and Women's Hospital, Boston, MA
- ⁴ Division of Geographic Medicine and Infectious Diseases, Tufts Medical Center, Boston, MA
- ⁵ Disease Control Bureau, Los Angeles County Department of Public Health, Los Angeles, CA
- ⁶ Department of Immunology and Infectious Diseases, Harvard T.H. Chan School of Public
- 14 Health, Boston, MA
- ⁷ Pandemic Sciences Institute, Nuffield Department of Medicine, University of Oxford, Oxford,
- 16 UK

17 Abstract

Epidemiological surveillance typically relies on reported incidence of cases or hospitalizations, 18 19 which can suffer significant reporting lags, biases and under-ascertainment. Here, we evaluated 20 the potential of viral loads measured by RT-qPCR cycle threshold (Ct) values to track epidemic trends. We used SARS-CoV-2 RT-qPCR results from hospital testing in Massachusetts, USA, 21 22 municipal testing in California, USA, and simulations to identify predictive models and covariates that maximize short-term epidemic trend prediction accuracy. We found SARS-CoV-2 Ct value 23 24 distributions correlated with epidemic growth rates under real-world conditions. We fitted gener-25 alized additive models to predict log growth rate or direction of reported SARS-CoV-2 case incidence using features of the time-varying population Ct distribution and assessed the models' abil-26 27 ity to track epidemic dynamics in rolling two-week windows. Observed Ct value distributions accurately predicted epidemic growth rates (growth rate RMSE ~ 0.039-0.052) and direction (AUC 28 29 $\sim 0.72-0.78$). Performance degraded during periods of rapidly changing growth rate. Predictive 30 models were robust to testing regimes and sample sizes; accounting for population immunity or symptom status yielded no substantial improvement. Trimming Ct value outliers improved perfor-31 32 mance. These results indicate that analysis of Ct values from routine PCR tests can help monitor 33 epidemic trends, complementing traditional incidence metrics.

34 Introduction

35 Epidemic monitoring and outbreak surveillance are vital public health functions, providing early warning of emerging threats, informing healthcare capacity planning and transmission control pol-36 37 icies, and helping to evaluate the effectiveness of interventions¹⁻⁴. A common approach to epidemic monitoring, exemplified during the COVID-19 pandemic, is to track the incidence of re-38 39 ported positive diagnostic tests, clinical cases^{5,6}, or deaths⁷. These data can inform key statistics such as the epidemic growth rate or effective reproductive number⁸⁻¹¹ and are fundamental to 40 41 nowcasting and forecasting an epidemic's trajectory¹²⁻¹⁴. However, these data streams can be substantially lagged, biased, and incomplete due to testing delays, capacity limitations, cost, and 42 changing test-seeking behavior^{15,16}. Thus, there has been growing interest in alternative data 43 sources, such as wastewater surveillance^{17,18}, internet search trends¹⁹, and digital contact trac-44 45 ing²⁰, that do not depend on large-scale testing of individuals,.

A novel data source for epidemic monitoring described during the COVID-19 pandemic is the 46 47 population-level distribution of viral loads among infected individuals, approximated using cycle threshold (Ct) values from reverse-transcription quantitative polymerase chain reaction (RT-48 49 qPCR) testing²¹⁻²⁴. For certain acute respiratory viruses such as SARS-CoV-2, a low Ct value 50 (high viral load) typically suggests that an individual was sampled early in their infection, whereas a high Ct value (low viral load) measurement suggests sampling later in infection^{25–27}. Thus, a 51 52 population-level sample of predominantly low Ct values (high viral loads) indicates that most sam-53 pled infections are of recent onset, corresponding to a growing epidemic, whereas a sample of predominantly high Ct values (low viral loads) corresponds to a declining epidemic consisting of 54 mostly late infections and post-infectious viral persistence²¹. Unlike count-based surveillance 55 56 methods, estimating epidemic growth rate based on the distribution of measured viral loads does 57 not depend on the number of positive tests.

58 Multiple studies have reported on the feasibility of using population-level Ct values to track SARS-59 CoV-2 epidemic trends.^{28–38} though it remains unclear under which conditions they are a practical 60 source of epidemiological information. While the relationship between sampled viral loads, viral 61 kinetics, and epidemic dynamics can be described mathematically under ideal conditions, in prac-62 tice there are several factors which complicate its application as a practical epidemic monitoring 63 tool. Measured Ct values are determined by a combination of biological factors, such as immunological history and variant causing the infection³⁹⁻⁴¹, and practical factors such as whether individ-64 65 uals are tested at a random point in their infection (e.g., asymptomatic screening) or around the time of peak viral load prompted by symptom onset⁴², demography of the tested population⁴³, 66 sample type⁴⁴, and RT-qPCR platform⁴⁵. Whether these factors are prohibitively confounding 67 when using Ct value distributions for epidemic monitoring has yet to be explored. 68

Here, we investigated the real-world feasibility of using SARS-CoV-2 Ct values to nowcast epidemic trajectories over three years of the COVID-19 pandemic. We first used synthetic datasets to benchmark nowcasting model performance and examined biological and logistical factors that might impede or improve nowcast accuracy. We then applied the same models to three real SARS-CoV-2 RT-qPCR testing datasets, collected across multiple geographic areas in the United States and under different population sampling strategies, to assess and inform the use of this approach in real-time estimation of epidemic growth rates.

76 **Results**

77 Correlation between epidemic growth rates and Ct value statistics using synthetic datasets

To understand how biological and practical factors might affect Ct-based nowcasting perfor-78 79 mance, we created several synthetic Ct value datasets using real population-level reported incidence curves for Massachusetts, USA, combined with a viral kinetics model parameterized by 80 longitudinal SARS-CoV-2 viral kinetics data, and sampling regimes representing a mixture of 81 symptom-driven testing and asymptomatic screening (see Methods). Using synthetic datasets in 82 83 this way allowed us to incorporate or exclude the effects of certain confounding factors on ob-84 served population-level Ct value distributions, in addition to the effect of the epidemic trajectory 85 itself (see Table S1).

86 We first simulated an ideal dataset assuming: 1) highly asymmetric viral kinetics, with a very short 87 growth phase and longer clearance phase; 2) low variation in observed viral load/Ct value for a given time-since-infection; and 3) a uniform probability of sampling an individual at any number 88 89 of days after infection or symptom onset. We varied each of these factors in turn, resulting in four alternative scenarios with either: 1) increased symmetry in viral kinetics, with a more similar 90 91 growth and clearance phase duration; 2) moderate variation in observed viral load/Ct value for a 92 given time-since-infection; 3) a low-variance, gamma-distributed delay between infection or symptom onset and sampling; and 4) a realistic baseline scenario combining all three factors (see 93 94 Supplementary Text 1).

All the synthetic datasets showed a clear negative correlation between the 7-day rolling average epidemic growth rate of cases and 7-day rolling average mean Ct value from the simulated symptomatic and asymptomatic samples, though the realistic baseline scenario showed the weakest correlation (**Figure S1**). Ct values from both symptom-based and random testing showed a relationship with epidemic growth rate (**Figure S1**), though Ct values observed through symptombased testing were typically lower and exhibited less variation.

101 With each synthetic dataset, we fit generalized additive models (GAM) with smoothing splines to 102 predict growth rates of cases as a non-linear function of daily mean and skewness of observed 103 Ct values. We also fit corresponding logistic GAMs to predict the epidemic direction, i.e., whether 104 incidence is growing or declining. We assessed in-sample fits of model-predicted vs. observed 105 growth rates and direction across the entire dataset, based on RMSE and AUC respectively. We then refit the models using separate training and testing subsets of the data. To approximate a 106 107 realistic application of the Ct-based approach in an ongoing epidemic, we fit the models using only the first 16 weeks of data and then performed rolling nowcasts with a two-week time horizon. 108 using the fitted model to estimate the epidemic growth rate and direction daily over the next two 109 weeks based on the Ct values reported during that time. At the end of each two-week window, we 110 re-fit the model using all Ct values and incidence data up to that time point, then nowcast the next 111 112 two-week window, and so on. As a sensitivity analysis, we compared RMSE and AUC with a fixed 113 train-test split date at the end of 2021 (Table S2).

With the ideal synthetic dataset, the GAMs closely tracked observed growth rates using Ct value
means and skew (Figure S2 & Figure S3; in-sample RMSE = 0.0191, approximately 10% of the
range in observed log incidence growth rates), as well as accurately predict epidemic direction
(in-sample AUC = 0.916). Nowcast accuracy over a rolling two-week window was slightly worse
than the in-sample predictive performance (mean across all nowcast windows, RMSE = 0.0206,
AUC = 0.867) but was still able to accurately track the epidemic over the full time period (Figure 1).

Model predictive performance was worse when using the realistic baseline synthetic dataset (**Figure 1, Figure S2 & Figure S3**; in-sample RMSE = 0.0319, AUC = 0.78; nowcast RMSE = 0.042, AUC = 0.698). The three factors examined individually had similar impacts on model performance; asymmetry of viral load trajectories caused the greatest increase in RMSE but only a slight decrease in AUC, while the distribution of delays between infection and sampling caused the

- smallest increase in RMSE but the largest decrease in AUC (Figure 1, Table 1). When these
- 127 models were applied to nowcasting growth rates in two-week increments, the greatest perfor-
- 128 mance reduction occurred when increasing the individual variation in Ct values for a given time-
- since-infection (**Table 1**).

Table 1. Predictive performance of GAMs using synthetic datasets, predicting per-day growthrates from daily Ct value statistics.

Dataset	RMSE, in-	RMSE, now-	AUC, in-	AUC, now-
	sample	cast	sample	cast
Ideal condition	0.0191	0.0206	0.916	0.867
Realistic kinetics	0.0245	0.0286	0.889	0.841
Realistic variation	0.024	0.0302	0.878	0.822
Realistic sampling	0.0237	0.027	0.865	0.824
Baseline condition	0.0319	0.042	0.78	0.698



133

Figure 1. Model-predicted (black) vs. observed (blue) log incidence growth rates for the five synthetic datasets, with model-predicted 95% confidence intervals (dark shading) and 95% prediction intervals (light shading). Predictions are from 2-week rolling nowcast, concatenated into a single time series. Vertical dashed line denotes the end of the initial training period. Variant era was included for comparability with later models; the synthetic datasets do not include any impact of viral variant on kinetics. WT=wild type.

139 Real-world relationship between observed Ct-value statistics and epidemic trajectories

Having established a baseline for model nowcasting performance using the synthetic data, we 140 141 next tested the nowcasting models on two RT-qPCR datasets: 1) routine hospital testing data from the Mass General Brigham hospital system in eastern Massachusetts (MGB), spanning Mar 2020-142 Feb 2023, and 2) municipal testing data from Los Angeles County, California (LAC), spanning 143 May 2020-Jul 2021 and Jan-Sep 2022. The MGB data came largely from mandatory screening 144 testing of outpatient, inpatients and emergency room admissions, while the LAC data were pri-145 marily symptom-driven voluntary testing (see Methods and Table S3). Both datasets contained 146 147 specimen collection dates and Ct values for SARS-CoV-2 positive results; LAC data also included vaccination status, symptom status, and symptom onset dates. MGB Ct values came from seven 148 platform/assay combinations, while Ct values from LAC data came from one PCR platform with 149 two possible assays (see Methods). 150

151 We limited our analysis to tests reporting Ct values, using the first available recorded Ct value for each infection episode (see Methods). The final analyzed sample included 104,534 (MGB) and 152 153 279,492 (LAC) Ct values. We also applied our method to a third, smaller set of testing data from the Tufts Medical Center, also in eastern Massachusetts, with a total sample of 10.214 Ct values. 154 155 We compared these Ct values against reported COVID-19 incidence for Massachusetts (MGB and Tufts) and Los Angeles County (LAC). We segment the data into four 'variant eras' based on 156 the SARS-CoV-2 variant known or believed to be dominant in the U.S. during different approxi-157 mate time periods, to allow for differences in viral kinetics by variant (see Methods). 158

159 Ct value distributions from both MGB and LAC datasets showed substantial variation over the 160 course of the pandemic (**Figure 2A & Figure 3A**). Reported COVID-19 incidence in both locations 161 varied over time as well (**Figure 2B & Figure 3B**), with large infection waves in the winters of 162 2020-21 and 2021-22, though the pattern of incidence was not synchronized across both settings.

163 While absolute incidence varied widely, incidence growth rates remained largely between ±0.2 164 throughout the course of the pandemic (**Figure 2B & Figure 3B**).

165 We found the mean and skewness of observed Ct value distributions (calculated daily over a seven-day moving window and excluding days with fewer than 10 Ct values reported) correlated 166 167 with the growth rate in reported incidence (Figure 2C & Figure 3C). Analysis of cross-correlation functions found Ct value distributions lagged incidence growth rate in the MGB data, with strong-168 est correlations at around 19-days lag (autocorrelation function, ACF= -0.462), and led incidence 169 growth rates for the LAC data, with strongest correlations at around 10-days lead (ACF = -0.062) 170 171 (Figure S4 & Figure S5). However, for real-time nowcasting, we focused on the relationship between same-day Ct values and incidence (i.e., lag=0 days; Figure 2C & Figure 3C), which still 172 showed high correlation. Higher incidence growth rates corresponded with lower same-day aver-173 174 age Ct values (Spearman's correlation coefficient: MGB Rho = -0.43, LAC Rho = -0.22) and with 175 positively skewed Ct distributions (MGB Rho = 0.35, LAC Rho = 0.43).



176

177 Figure 2. Ct values from the Mass General Brigham hospital system and corresponding reported COVID-19 incidence in Massachusetts, USA. (A) Weekly Ct value guantiles over time, showing 178 weekly median Ct value and 50/80/90/95% guantiles. (B) 7-day rolling average reported incidence 179 (grey bars), growth rate in 7-day rolling average reported incidence (grey line), and smoothed 180 growth rate (blue line). Background is shaded by time periods of different variant dominance. 181 Vertical dashed line demarcates the test-train split. (C) Incidence growth rate compared to 182 smoothed daily mean and skewness of Ct value distributions. Colored lines and shaded grey 183 regions show fitted cubic spline GAMs with 95% confidence intervals, stratified by period of variant 184 185 dominance.



187 Figure 3. Ct values from Los Angeles County and corresponding reported COVID-19 incidence. (A) Weekly Ct value quantiles over time, showing weekly median Ct value and 50/80/90/95% 188 quantiles. (B) 7-day rolling average reported incidence (grey bars), growth rate in 7-day rolling 189 average reported incidence (grey line), and smoothed growth rate (blue line). Background is 190 shaded by time periods of different variant dominance. Vertical dashed line demarcates the test-191 train split. (C) Incidence growth rate compared to smoothed daily mean and skewness of Ct value 192 distributions. Colored lines and shaded grey regions show fitted cubic spline GAMs with 95% 193 194 confidence intervals, stratified by period of variant dominance.

Nowcasting epidemic growth rates using Ct values in Massachusetts, USA and Los Angeles
 County, USA

We next re-trained the same GAM models used with synthetic data to the MGB and LAC dataset using smooth functions of mean and skewness of Ct values to predict log incidence growth rates, with corresponding logistic models to predict epidemic direction. Model predictions were compared against observed values first in-sample across the entire dataset then over a rolling twoweek nowcast window, as well as with a single fixed train-test split date at the end of 2021.

In both datasets, this simple model achieved in-sample prediction accuracy for incidence growth rate only slightly worse than performance on the realistic synthetic data, with relatively small absolute errors (MGB RMSE = 0.0451; LAC RMSE = 0.0335, see **Figure S6-S9**, **Table S4**, and **Table 2**). Corresponding logistic regression models successfully discriminated growing from declining incidence (Area under the curve: MGB AUC = 0.785, LAC AUC = 0.843).

207 The models were able to nowcast growth rates, in two-week increments with models periodically refitted to more recent data, with accuracy slightly worse than in-sample model fits (MGB RMSE 208 = 0.0523, LAC RMSE = 0.039) (Figure 4A & Figure 5A). This level of nowcast accuracy was 209 210 likewise only slightly worse than nowcasting performance with realistic synthetic data. While av-211 erage prediction error was relatively small, comparable to in-sample model error and to prediction 212 error with realistic synthetic data, accuracy was highly variable from one two-week window to the next (Figure 4B & Figure 5B). Nowcast accuracy was comparable to model performance over a 213 214 fixed multi-month prediction window, slightly better for one dataset and worse for the other (MGB 215 RMSE = 0.047, LAC RMSE = 0.0458) (Table S2). Nowcast predictions of epidemic direction were 216 slightly worse than in-sample ones (MGB AUC = 0.723, LAC AUC = 0.784) and outperformed the directional discrimination test with realistic synthetic data. In addition, over all two-week nowcast 217 windows combined, model-predicted growth rates correlated moderately well with observed ones 218 (Spearman's Rho: MGB Rho = 0.398, LAC Rho = 0.556). 219

Table 2. Predictive performance of the selected GAM using data from MGB and LAC, predicting

221 per-day growth rates from daily Ct value statistics.

Dataset RMSE

AUC

	In-sample	Nowcast	Periods of rapid change in growth rate	In-sample	Nowcast	Periods of rapid change in growth rate
MGB	0.0451	0.0523	0.0645	0.785	0.723	0.722
LAC	0.0335	0.039	0.0471	0.843	0.784	0.772



Figure 4. (A) Model-predicted (black) vs. observed (blue) log incidence growth rates for MGB data, with 95% confidence intervals (dark shading) and 95% prediction intervals (light shading). (B) RMSE of predicted vs. observed log incidence growth rates for each 2week nowcasting window. "Inflection periods" refer to times when the absolute smoothed log incidence growth rate exceeded 0.025 over a one-week period, marked with points above each subplot.



Figure 5. (A) Model-predicted vs. observed log incidence growth rates and RMSEs for LAC data. Model-predicted (black) vs. observed (blue) log incidence growth rates for MGB data, with 95% confidence intervals (dark shading) and 95% prediction intervals (light shading). **(B)** RMSE of predicted vs. observed log incidence growth rates for each 2-week nowcasting window. "Inflection periods" refer to times when the absolute smoothed log incidence growth rate exceeded 0.025 over a one-week period, marked with points above each subplot.

Nowcasting performance during time periods of rapid change in growth rate

To assess nowcasting performance during periods of rapid change in the epidemic trajectory, we identified times when the absolute smoothed incidence growth rate exceeded 0.025 over a oneweek period. This definition captured 30.1% of nowcast dates for the MGB data (284/944 days) and 17.5% for LAC (98/560 days). We then recalculated in-sample and out-of-sample prediction accuracy for growth rate and epidemic direction during just these periods.

Across both datasets, prediction error over periods of rapid change was greater than over the whole nowcast period (MGB RMSE = 0.0645 [change] vs. 0.0523 [nowcast], LAC RMSE = 0.0471 [change] vs. 0.039 [nowcast]; see **Table 2**). However, directional prediction accuracy was comparable between periods of rapid change and the whole nowcast period (MGB AUC = 0.722 vs. 0.723, LAC AUC = 0.772 vs. 0.784).

Nowcasting performance with variable sample size and outlier removal

We assessed the sensitivity of nowcasting performance to sample size both by randomly 246 247 downsampling the MGB dataset (100 random draws) and by analyzing a third, smaller dataset from Tufts Medical Center using the same response variable (i.e., log incidence growth rates for 248 249 Massachusetts) but with approximately 10% of the total sample size of the MGB data (see Meth-250 ods; Figure S10, S11). In most cases, prediction accuracy for incidence growth rate was comparable with the downsampled datasets and the equivalent full datasets (Figure 6; see also Table 251 252 **S5**). Only with 10% of the full dataset (but not with the Tufts dataset) did nowcasting accuracy 253 degrade appreciably; with 50-75% downsampling or a daily maximum of 25 positive samples. accuracy improved compared to baseline. Likewise, directional prediction accuracy was generally 254 255 similar between downsampled and full datasets, with substantially worse accuracy only for the 256 10% downsample. Improved accuracy may reflect reduced influence of outliers – downsampling 257 the full dataset tends to exclude the days with smallest sample sizes, which are otherwise given equal weight in model training to days with more observations, while sub-sampling each day's 258

observations reduces the impact of outliers on each day's observed Ct value distribution. To test
this, we examined model performance with trimming of outlier Ct values from each day's observed
data. Trimming outliers reduced prediction error with 2.5%, 5%, and 10% trims (Figure 6, Table
S5), while 2.5% and 5% trims also improved directional prediction accuracy.



263

Figure 6. Model performance for downsampled MGB and full Tufts datasets. Baseline comparison metrics are re-calculated for only the days included in each downsampled dataset's nowcast. For proportional and daily max downsampling, both downsampled and baseline performance are averaged over 100 random draws (and their corresponding days included). Trim percentages indicate quantiles trimmed from each end of daily Ct value distributions (i.e., 5% trim yields the 5-95 percentile range of Ct values).

270 Sensitivity analysis

271 The reason for testing (e.g., symptom driven testing vs. screening asymptomatic outpatients) is 272 expected to result in different distributions of observed Ct values due to variation in when individ-273 uals are tested during their infection; therefore, the relationship between Ct values and epidemic 274 growth rate is expected to differ correspondingly. In addition, in the MGB data, individuals were swabbed differently and tested on different PCR platforms depending on their reason for seeking 275 276 healthcare, including a mixture of patients tested as outpatients, inpatients and in the emergency 277 room. To understand the impact of these factors on the modelled relationship between Ct values 278 and growth rate, we assessed performance of GAMs using only 1) MGB data from routine outpa-279 tient screening, the majority of whom were sampled in the same way and tested on the same PCR platform (Figure S12); 2) LAC data stratified by symptom status (symptomatic vs. asympto-280 281 matic vs. no known symptom status); 3) LAC data from tests conducted on asymptomatic individ-282 uals and those without known symptom status; and 4) LAC data from unvaccinated individuals with no known previous SARS-CoV-2 infection. In all cases, we compared performance to the 283 base model for the respective data source. The relationship between Ct values and growth rate 284 285 appeared to differ when subsetting or stratifying by these variables (Figure S13), but including these stratifications in the model did not always improve predictive performance. Restricting to 286 outpatient tests only improved prediction error compared to baseline (nowcast RMSE = 0.0494 287 vs. 0.0523 base), whereas incorporating symptom status or immune history slightly worsened 288 prediction error (nowcast RMSE = 0.454 for symptom-stratified, 0.0415 for asymptomatic/no 289 symptom status only, 0.0401 for immunologically naïve only, vs. 0.039 base, see Table S6). 290

291 Discussion

Under real-world conditions, simple generalized additive models using the mean and skewness
of recorded Ct values could nowcast (log) incidence growth rates with prediction errors (RMSE)
of approximately 0.04-0.05. Across both settings (Massachusetts and Los Angeles County),

growth rates generally varied between approximately ± 0.2 , so this level of accuracy in modelled estimates, while not highly precise, is nonetheless informative. These models are also able to identify if incidence is growing or shrinking with AUC greater than 0.7, substantially better than chance.

Nowcast accuracy over two-week time horizons is slightly worse than the quality of in-sample model fits, especially early into the emergence of new dominant viral variants whose effect cannot yet be accurately estimated. During periods of rapid change in incidence growth rate (e.g., just as a new outbreak wave is developing), nowcast accuracy for growth rate is slightly worse, possibly due to larger absolute growth rates during such periods. Crucially, however, directional predictions remain moderately accurate during those times.

305 Our results support the theoretical expectation that epidemic dynamics influence population-level 306 viral load distributions, and therefore can be inferred from them²¹. They also corroborate the find-307 ings from other settings, where Ct values have been used successfully to infer epidemic growth rates or reproduction numbers^{28–38}. Our analysis builds on these studies with one of the largest 308 empirical tests of this nowcasting approach to date using data from two locations in the USA over 309 310 a three-year period. Epidemic growth rates and directions were accurately nowcasted using both 311 datasets, despite showing different Ct value trends and capturing different populations, highlight-312 ing the generalizability of this approach. Furthermore, these data covered a long-time window and included periods of different variant dominance and population immunity, suggesting Ct values 313 314 could continue to augment infectious disease surveillance as SARS-CoV-2 epidemiology contin-315 ues to change.

In practice, several factors can confound the relationship between Ct values and epidemic dynamics (measured here as growth rate of case incidence), including testing delays, sampling regimes (i.e., community-based random testing vs. testing patients in hospital), symptomatic (diagnostic) vs. asymptomatic (screening) testing, immunological history, and the inherent individual-

320 level variability in SARS-CoV-2 viral kinetics. Our synthetic data analyses help disambiguate these confounding factors by comparing degradation in predictive performance between different 321 322 synthetic datasets. Predictive performance was slightly worse in the real datasets compared to a 323 'realistic' synthetic dataset. One key contributor is that Ct values from the real datasets were col-324 lected using multiple RT-qPCR assays and/or platforms and were not standardized and may gen-325 erate different Ct values for the same underlying viral load, limiting the comparison of Ct values 326 across platforms and assays (see Methods)^{45–47}. Additionally, the data-generating model for our 327 'realistic' synthetic dataset did not incorporate the impact of vaccination or past infection which affect individual viral load trajectories^{40,48}, potentially contributing to the differences in performance 328 between models with empirical vs. synthetic data. 329

330 Our synthetic data analysis also highlights the importance of considering the delay between in-331 fection and sampling an individual in determining the population-distribution of Ct values. Funda-332 mentally, the relationship between population-level epidemic dynamics and viral load distributions arises because individuals' viral loads reflect times since infection²¹, and hence cross-sectional 333 334 distributions of viral loads (or Ct values) reflect the distribution of times-since-infection among currently infected individuals, similar to the relationship between incidence and prevalence. This 335 336 relationship can be readily described mathematically if individuals are randomly sampled, with a uniform probability of sampling any time after infection. Random cross-sectional samples captur-337 ing infections at random points in their infection are rare (see ^{22,31,36,49}) but are reasonably well 338 approximated in our datasets by routine screening of hospital outpatients. However, a more real-339 340 istic sampling delay distribution – such as if individuals tend to be tested shortly after suspected exposure or developing symptoms – biases the probability of sampling over time since infection 341 and dilutes the signal of infection age. Symptom-driven testing where individuals are tested due 342 to recent symptom onset beginning at around the same time as peak viral load, is the most com-343 344 mon source of data used for epidemiological surveillance, which reduces any epidemic signal in

the population-level Ct distribution. In the extreme, if individuals were sampled with the same delay following infection, then any observed variation in viral loads would arise from random individual variation at a single time-since-infection rather than reflecting a distribution of times-sinceinfection among current infections. Changes in public health recommendations around testing and screening algorithms, such as recommendations around pre-travel testing or hospital admissions screening, may therefore change the relationship between population Ct values and epidemic dynamics, which may bias Ct-based epidemiological estimates if not accounted for.

PCR platform differences and nonrandom sampling regimes are both addressable challenges, at 352 353 least in principle. Ct value data could be calibrated across platforms and assays using standardized samples. Random surveillance sampling could reduce the bias in testing delay found with 354 355 symptom-driven testing. True random sampling may be important, as voluntary testing by asymp-356 tomatic individuals may still show some bias in testing delays (**Table S6**). When we approximated 357 these changes by subsetting one of our datasets to only results from outpatient screening tests, which were largely collected and analyzed the same way (Figure S12), we found small improve-358 359 ments in model predictive performance compared to using the full, mixed dataset (Table S6). While random surveillance sampling at low prevalence may yield very few infections detected, 360 361 nowcasting accuracy was not severely degraded even with substantially reduced sample sizes (Figure 6). Both these changes would improve the accuracy of simple Ct-based nowcasting mod-362 els. Even absent such logistical solutions, however, we found the simple statistical heuristic of 363 364 trimming outliers (2.5-5%) from daily observed Ct values improves nowcasting accuracy (Figure 365 6).

Beyond confounding factors, it is plausible that growth rate of reported COVID-19 cases may not be the most accurate benchmark against which to compare Ct value distributions. First, symptomatic cases occur and are reported with a lag relative to infections, and may be affected by changes in testing behavior, for example with the increased availability of home-based rapid

antigen tests. Alternative benchmarks, such as growth rate in hospitalizations, mortality, or wastewater viral loads, may therefore yield stronger relationships (possibly with some time-shifting); investigating these relationships would be a fruitful avenue for further research. In addition, geographically aggregated incidence may mask heterogenous outbreak trajectories at finer scale, e.g., city or even neighborhood level. Such finer-scale incidence data may yield cleaner relationships with Ct value distributions, especially if matched to the catchment areas for the Ct value data collection process.

377 Another challenge for modeling Ct value dynamics is the choice of mathematical model to capture 378 the relationship between observed Ct values and underlying epidemic growth rates. The link between epidemic dynamics and viral loads observed through random cross-sectional surveillance 379 380 can be described precisely based on the convolution of the infection incidence curve and viral kinetics curve^{21,31,36}. In contrast, viral loads observed through non-random or convenience sam-381 382 ples, such as symptom-driven testing, arise from complex data generating processes which are difficult to describe mathematically, and thus past studies, including ours, tend to favor regression 383 models to estimate epidemic dynamics from observed Ct values^{30,33,34}. Future work should focus 384 on more complex statistical methods that take into account the time-series nature of the data³⁷, 385 386 the non-linear and potentially non-monotonic relationship between Ct values and growth rates, 387 and combine multiple data streams to provide more accurate predictions of epidemic dynamics³⁶.

Tracking epidemic growth rates in near-real-time remains an important challenge for public health surveillance. Our analyses show that simple Ct-based models can accurately track SARS-CoV-2 epidemic growth rates, highlighting their potential use in augmenting infectious disease surveillance systems. Ultimately, their greatest strength lies in their speed and simplicity. The models presented here are conceptually straightforward and computationally lightweight, easy to implement even in resource-constrained settings, and, unlike wastewater testing, are reliant only on data already routinely collected as part of screening or diagnostic testing. Our analyses show that

395	they retain their accuracy even with limited sample sizes or during periods of rapid change in
396	epidemic trajectories, such as during the transition from the end of one epidemic wave to the start
397	of the next one, and so could provide rapid situational awareness as outbreak waves emerge.
398	Further research could examine how Ct-based estimates of epidemic trajectories complement
399	other, orthogonal indicators such as wastewater surveillance, as well as potential applications to
400	different viral pathogens with well-characterized viral kinetics such as influenza or RSV ^{50,51} .

402 Methods

403 Study settings & data sources

404 *Massachusetts*

405 Massachusetts Ct value data comes primarily from testing in 16 hospitals in the Mass General Brigham hospital system, with a catchment area largely in eastern Massachusetts. The full da-406 407 taset comprises 2,671,041 SARS-CoV-2 test results, with specimen collection dates ranging from 3 Mar 2020 to 23 Feb 2023, of which 161,273 were positive. There were 3531 individuals who 408 409 appeared to experience repeat infections (defined as >60 days between positive results), of which 72 individuals had 2 or more repeat infections. As we could not rule out long COVID or other 410 idiosyncratic viral kinetics, we drop these 72 individuals from the final dataset. Limiting to results 411 412 reporting Ct values and first reported Ct values for each confirmed case yields the final sample of 413 104,534 Ct values used in this analysis (**Table S3**), of which the earliest specimens were collected on 31 Mar 2020. 414

Samples are from a combination of routine outpatient (77,700; 74.3% of samples) and inpatient (7,311; 7.0%) screening and diagnostic tests, as well as ER patient testing (19,523; 18.7% of samples); while not entirely random nor representative, routine screening tests suffer less selfselection bias than symptom-based or voluntary testing. We did not have access to information on patients' vaccination or infection history, infecting variant, or symptom status.

The final sample includes specimens collected from nasal and nasopharyngeal swabs (approx. 2:1 ratio). Specimens were processed using seven different RT-qPCR platform/assay combinations (**Table S3**), variously targeting E/N/N1/N2/ORF1ab genes. For the main analysis here, Ct values were pooled across platforms/assays; where a single result reported Ct values for multiple target genes, the lowest value was used.

Daily confirmed case counts for Massachusetts were obtained from the Massachusetts Depart ment of Public Health COVID-19 dashboard⁵².

We also analyzed a secondary dataset of Ct values from Tufts Medical Center in Boston, Massachusetts for comparison. This dataset comprised 84,848 test results with collection dates ranging from 18 Feb 2021 to 31 Oct 2022, of which 10,338 were positive. Filtering the reported test results using the same criteria as used for the MGB data yielded a final sample of 10,214 Ct values used here. **Figure S11** summarizes the reported Ct value distributions over time and compares these to reported COVID-19 incidence.

433 Los Angeles County

LAC Ct value data comes from municipal COVID-19 testing sites operated by the LAC Department of Public Health and Department of Health Services, comprising approximately 10% of all municipal testing conducted in LAC during the sample period. The full dataset comprises 330,034 SARS-CoV-2 positive test results, with specimens collected over two time periods – 21 May 2020 to 27 Jul 2021, and 30 Dec 2021 to 29 Sep 2022. (Note: data were unavailable for the intervening period.) The data contain an infection episode identifier; limiting to the first reported Ct value for each infection episode yields the final sample of 279,492 Ct values used in this analysis.

441 The final sample includes specimens collected through nasal, nasopharyngeal, and oral swabs, 442 and analyzed by Fulgent Genetics using an in-house platform and ThermoFisher QuantStudio™ 6 and 7 PCR systems. Two RT-qPCR assays were used; before mid-Nov 2020, analyses used 443 exclusively LOINC 94531-1 targeting N1 and N2 genes, while subsequently the majority of anal-444 445 yses used LOINC 94533-7 targeting the N gene. Where a single result reported Ct values for multiple target genes, the lowest value was used. Symptom status was reported for approximately 446 447 75% of the sample, of which in turn approximately 75% (56% of the full sample) are reported as symptomatic for COVID-19 (Table S3). For symptomatic cases, most specimens were collected 448

1-10 days after symptom onset (modal delay of 3 days). The sample also included vaccination
status, with approximately 24% of results coming from vaccinated (partially, fully, or boosted) individuals (**Table S3**).

452 Daily confirmed case counts were obtained from the LAC DPH COVID-19 dashboard⁵³.

453 Synthetic datasets

454 We built on a previously published model to simulate realistic Ct value distributions that would be expected under testing and sampling schemes similar to real-world data²¹. Full details of the sim-455 456 ulation framework are given in Supplementary Text 1. First, we parameterized a viral kinetics 457 model describing the expectation and distribution of Ct values over all days following infection 458 using previously published longitudinal SARS-CoV-2 testing data (Figure S15, Table S7)⁴⁸. This is a piecewise linear model governed by a set of control points determining the time from infection 459 460 to peak viral load, time from peak viral load to an inflection point at a high Ct value, and a longer-461 term clearance rate with a daily probability of full clearance. Second, we simulated approximately 462 2 million infections with infection times distributed based on the reported incidence of COVID-19 cases in Massachusetts between 5 March 2020 and 25 Feb 2023. Third, we simulated a surveil-463 lance system as a mixture of random testing (i.e., symptom-independent) and symptom-based 464 465 testing (individuals are tested with a random delay following a randomly generated incubation 466 period). Combining these three simulation steps gave a synthetic dataset of Ct values for a mixture of asymptomatic and symptomatic individuals tested at various times post infection and over 467 a multi-wave SARS-CoV-2 epidemic (Figure S1). Different scenarios were captured by changing 468 the parameters used either for the viral kinetics model or sampling delay distribution (Figure S16). 469

470 Statistical methods

We calculated daily incidence-based growth rates as the natural log-transformed ratio of 7-day
moving average new reported cases for each day to the 7-day moving average for the preceding
day:

474
$$y_t = \ln \frac{\sum_{k=0}^{6} f_{(t-k)}}{\sum_{k=1}^{7} f_{(t-k)}}$$

475 where y_t is incidence growth rate and f_t is daily incidence at time t. We defined epidemic direction

476 as growing when $y_t > 0$ and declining when $y_t \le 0$.

477 Classifying time periods of rapid incidence change

To identify periods of rapid change in incidence growth rate, we first smoothed the daily incidence growth rate (as defined above) using a centered 7-day moving average:

480
$$y'_t = \frac{1}{7} \sum_{k=-3}^{3} y_{(t+k)}$$

We then identified times when the absolute change in smoothed log incidence growth rate y'_t equals or exceeds 0.025 over a one-week period, denoting the midpoint days of those weeks as periods of rapid change. That is, time *t* is defined as having rapid change in incidence if and only if $|y'_{(t+3)} - y'_{(t-3)}| \ge 0.025$.

485 Growth rate & epidemic direction models

486 We modeled incidence growth rate using a generalized additive model (GAM) incorporating the 487 mean and skewness of Ct values:

488
$$\ln y_t = \beta_0 + s_{\bar{x}}(\bar{x}_t) + s_g(g_t) + \beta_v v_t$$

Where $s_{\bar{x}}$ and s_g are smoothing functions fitted using cubic regression splines⁵⁴, and \bar{x}_t and g_t are the 7-day rolling averages at time *t* of the daily mean and skewness respectively of Ct values

491 from samples collected or over the window from time t to t - 6, excluding days with fewer than 10 Ct values reported. v_t is a categorical variable identifying the SARS-CoV-2 variant known or 492 believed to be dominant in the U.S. during different approximate time periods. For our datasets, 493 we designated four such variants / time periods: wild type (up to 30 Nov 2020). Alpha (01 Dec 494 495 2020 to 31 May 2020), Delta (01 Jun 2020 to 03 Dec 2021), and Omicron (04 Dec 2021 onwards). 496 We used this rough approximation rather than relying on more direct and detailed observations, e.g. sequencing data linked to our datasets, to better represent a realistic use case for the Ct-497 498 based method such as a small municipal public health department. In such cases, resources for extensive sequencing may not be available, necessitating reliance on broader national trends. 499 500 When encountering new variant[s] in a nowcasting or testing period not present in training data, 501 our models use a realistic decision rule of making predictions based on the last known variant 502 from training data.

503 We model epidemic direction using logistic regression models equivalent to the GAMs used for 504 incidence growth rate.

505 To determine our choice of model, we tested a series of log-linear regression models and GAMs, 506 using different predictors (daily Ct mean, standard deviation, and skewness), functional forms 507 (log-linear vs. cubic regression splines), and variant interaction terms. We fitted these models to 508 the baseline synthetic dataset and compared their AIC as well as in-sample and nowcasting per-509 formance (see below). There was a clear bias-variance tradeoff between models; more flexible 510 model specifications yielded better AIC and in-sample fit, at the cost of worse out-of-sample or nowcasting performance (see Table S2 and Figure S14). We ultimately selected the final model 511 512 using mean and skewness with a cubic spline, as the theoretical relationship between cross-513 sectional Ct values and epidemic growth rates is non-linear and depends on the *distribution* of Ct values observed; short of fitting the growth rate model to the entire distribution of observed values, 514

using mean and skewness provides a parsimonious way to include information about the shapeof the distribution in the model.

517 Evaluating model performance

To evaluate the performance of Ct-based nowcasting models, we conducted two model validation tests. First, we fitted the main models to each dataset using only data up to 31 Dec 2021, then used the fitted model to predict incidence growth rates and epidemic direction for the remainder of each dataset (from 01 Jan 2022 onwards), based on observed Ct values reported in each dataset. We assessed prediction performance using RMSE between model-predicted and observed incidence growth rates, as well as AUC for directional predictions from the logistic regression model.

525 Next, we conducted a 'rolling' nowcast test, intended to simulate a realistic application of this approach. For each dataset, we trained the main models on the first 16 weeks of available data, 526 using the models thus fitted to predict incidence growth rates and epidemic direction over the 527 528 following 2-week period using only reported Ct value statistics. We then re-fit the models incorporating those two weeks of incidence data (i.e., up to 18 weeks) and predict the subsequent 2-529 530 week period, repeating this re-fitting and prediction procedure in 2-week increments up to the end 531 of each dataset. We report prediction performance as RMSE or AUC across all 2-week prediction periods concatenated into a single prediction time series for each dataset and model, while de-532 533 tailed period-by-period performance is reported in the online repository at ⁵³.

534 Impact of reduced sample size and outliers on Ct-based growth rate estimation

As a sensitivity analysis, we repeated the rolling nowcast analyses using artificially down-sampled datasets. We generated downsampled versions of the dataset in two ways: 1) by randomly drawing 10/25/50/75% of the total test results available, or 2) by limiting the maximum number of positive test results for each day to 25/50/100, discarding any additional tests. We then reassessed

539 nowcasting performance on each of these downsampled datasets. We repeated this analysis with 540 100 different randomly downsampled datasets for each size, taking the mean of model perfor-541 mance metrics over the 100 draws at each size. We also compared nowcasting performance with 542 a similar analysis using a third, smaller dataset from Tufts Medical Center, which uses the same 543 response variable data as the MGB dataset (i.e., log incidence growth rates for Massachusetts) but has approximately 10% of the total sample size. The downsampling process can result in 544 some days being excluded from the downsampled dataset model's nowcast. Nowcasting perfor-545 546 mance can vary considerably from day to day, with outlier days having disproportionate impact. To ensure fair comparison of the impact of downsampling on model accuracy, rather than the 547 impact of certain days being excluded as an indirect result of the downsampling process, we 548 recalculated performance metrics for the baseline model's nowcasts based on just the days in-549 cluded in any given downsampled model's nowcasts, once again taking the mean of model per-550 551 formance metrics over the 100 different baseline subsets included for each sample size.

To assess the impact of outliers on nowcasting performance, we trimmed daily observed Ct value distributions by 2.5/5/10% (yielding 95/90/80% ranges) before calculating Ct value distribution statistics, using the trimmed data for both training and nowcasting. Repeat draws were not required as the trimming is deterministic. As with the downsampling analysis, we recalculated baseline model performance metrics for only days included at each trim level.

557 Data & code availability

558 Data and analysis code are available online at https://github.com/gradlab/ct-nowcasting [NOTE: 559 we will update this to a Zenodo DOI before publication].

560 Acknowledgements & financial disclosures

JAH is supported by a Wellcome Trust Early Career Award (grant 225001/Z/22/Z). This work was supported in part by the Francis P. Tally, MD, Fellowship in the Division of Geographic Medicine and Infectious Disease (JAP). This project has been funded in part by contract 200-2016-91779

with the Centers for Disease Control and Prevention (CDC). Disclaimer: The findings, conclusions, and views expressed are those of the authors and do not necessarily represent the official position of the CDC. The authors also thank Jason Cheng and Hanlin (Harry) Gao of Fulgent Genetics for assistance with data for the analysis.

- 568 All authors declare no competing interests. No authors nor our institutions received any payments
- or services in the past 36 months from a third party that could be perceived to influence, or give
- 570 the appearance of potentially influencing, the submitted work.
- 571 Ethics guidelines

572 The authors declare that all relevant ethical guidelines have been followed and all necessary IRB 573 and/or ethics committee approvals have been obtained.

574 Author contributions

575 TYL, YHG, and JAH conceptualized the project. TYL and JAH designed the analyses, developed

the code, and created the visualizations. TYL, SK, JP, MH, THK, and PD prepared data. SK, SD,

577 RF, and YHG provided resources and contributed to analysis design and interpretation. YHG pro-

vided primary supervision and funding support. TYL and JAH wrote the first draft. All authors

579 provided critical review and revision of the text and approved the final version.

580 **References**

- Bhatia, S. *et al.* Lessons from COVID-19 for rescalable data collection. *Lancet Infect. Dis.* e383–e388 (2023).
- 583 2. Cori, A. et al. Key data for outbreak evaluation: building on the Ebola experience. Philos.
- 584 Trans. R. Soc. Lond. B Biol. Sci. **372**, 20160371 (2017).
- Lipsitch, M. *et al.* Infectious disease surveillance needs for the United States: lessons from
 COVID-19. *arXiv* [*cs.CY*] (2023).
- 4. Lipsitch, M. *et al.* Improving the evidence base for decision making during a pandemic: the example of 2009 influenza A/H1N1. *Biosecur. Bioterror.* **9**, 89–115 (2011).
- 589 5. https://ourworldindata.org/covid-cases.
- 590 6. UKHSA data dashboard. https://ukhsa-dashboard.data.gov.uk/.
- Flaxman, S. *et al.* Estimating the effects of non-pharmaceutical interventions on COVID-19
 in Europe. *Nature* 584, 257–261 (2020).
- Solution Series Serie
- Parag, K. V., Thompson, R. N. & Donnelly, C. A. Are epidemic growth rates more informative
 than reproduction numbers? *J. R. Stat. Soc. Ser. A Stat. Soc.* (2022) doi:10.1111/rssa.12867.
- Abbott, S. *et al.* Estimating the time-varying reproduction number of SARS-CoV-2 using
 national and subnational case counts. *Wellcome Open Res.* 5, 112 (2020).
- Charniga, K. *et al.* Updating reproduction number estimates for mpox in the Democratic
 Republic of Congo using surveillance data. *Am. J. Trop. Med. Hyg.* **110**, 561–568 (2024).
- 12. Charniga, K. *et al.* Nowcasting and forecasting the 2022 U.S. mpox outbreak: Support for
- 603 public health decision making and lessons learned. *Epidemics* **47**, 100755 (2024).

- 604 13. Günther, F., Bender, A., Katz, K., Küchenhoff, H. & Höhle, M. Nowcasting the COVID-19
 605 pandemic in Bavaria. *Biom. J.* 63, 490–502 (2021).
- Reich, N. G. *et al.* A collaborative multiyear, multimodel assessment of seasonal influenza
 forecasting in the United States. *Proc. Natl. Acad. Sci. U. S. A.* **116**, 3146–3154 (2019).
- 15. Rahmandad, H., Lim, T. Y. & Sterman, J. Behavioral dynamics of COVID-19: estimating
 underreporting, multiple waves, and adherence fatigue across 92 nations. *Syst. Dyn. Rev.*37, 5–31 (2021).
- 16. Tsang, T. K. *et al.* Effect of changing case definitions for COVID-19 on the epidemic curve
- and transmission parameters in mainland China: a modelling study. *Lancet Public Health* 5,
 e289–e296 (2020).
- Huisman, J. S. *et al.* Wastewater-based estimation of the effective reproductive number of
 SARS-CoV-2. *Environ. Health Perspect.* **130**, 57011 (2022).
- 18. Tisza, M. J. *et al.* Virome sequencing identifies H5N1 avian influenza in wastewater from nine
 cities. *bioRxiv* 2024.05.10.24307179 (2024) doi:10.1101/2024.05.10.24307179.
- Stolerman, L. M. *et al.* Using digital traces to build prospective and real-time county-level
 early warning systems to anticipate COVID-19 outbreaks in the United States. *Sci. Adv.* 9,
 (2023).
- Kendall, M. *et al.* Drivers of epidemic dynamics in real time from daily digital COVID-19
 measurements. *Science* 385, (2024).
- 623 21. Hay, J. A. *et al.* Estimating epidemiologic dynamics from cross-sectional viral load
 624 distributions. *Science* 373, eabh0635 (2021).
- Walker, A. S. *et al.* CT threshold values, a proxy for viral load in community sars-cov-2 cases,
 demonstrate wide variation across populations and over time. *Elife* **10**, (2021).
- 23. Penney, J., Jung, A., Koethe, B. & Doron, S. Cycle threshold values and SARS-CoV-2:
 Relationship to demographic characteristics and disease severity. *J. Med. Virol.* 94, 3978–
- 629 **3981 (2022)**.

- 630 24. Sala, E. *et al.* Systematic review on the correlation between SARS-CoV-2 real-time PCR
- 631 cycle threshold values and epidemiological trends. *Infect. Dis. Ther.* **12**, 749–775 (2023).
- 632 25. Kissler, S. M. *et al.* Viral dynamics of acute SARS-CoV-2 infection and applications to
 633 diagnostic and public health strategies. *PLoS Biol.* **19**, e3001333 (2021).
- 634 26. Tsang, T. K. *et al.* Influenza A Virus Shedding and Infectivity in Households. *J. Infect. Dis.*635 **212**, 1420–1428 (2015).
- Brint, M. E. *et al.* Prolonged viral replication and longitudinal viral dynamic differences among
 respiratory syncytial virus infected infants. *Pediatr. Res.* 82, 872–880 (2017).
- 28. Harrison, R. E. *et al.* Cycle Threshold Values as Indication of Increasing SARS-CoV-2 New
 Variants, England, 2020-2022. *Emerg. Infect. Dis.* **29**, 2024–2031 (2023).
- 640 29. Musalkova, D. *et al.* Trends in SARS-CoV-2 cycle threshold values in the Czech Republic
 641 from April 2020 to April 2022. *Sci. Rep.* **13**, 6156 (2023).
- 642 30. Lin, Y. *et al.* Incorporating temporal distribution of population-level viral load enables real-time
 643 estimation of COVID-19 transmission. *Nat. Commun.* **13**, (2022).
- Aguilar Ticona, J. P. *et al.* Extensive transmission of SARS-CoV-2 BQ.1* variant in a
 population with high levels of hybrid immunity: A prevalence survey. *Int. J. Infect. Dis.* 139,
 159–167 (2024).
- 32. Andriamandimby, S. F. *et al.* Cross-sectional cycle threshold values reflect epidemic
 dynamics of COVID-19 in Madagascar. *Epidemics* 38, 100533 (2022).
- 33. Alizon, S. *et al.* Epidemiological and clinical insights from SARS-CoV-2 RT-PCR crossing
 threshold values, France, January to November 2020. *Euro Surveill.* 27, (2022).
- 34. Yin, N. *et al.* Leveraging of SARS-CoV-2 PCR cycle thresholds values to forecast COVID-19
 trends. *Front. Med. (Lausanne)* 8, (2021).
- 35. Khalil, A. *et al.* Weekly Nowcasting of New COVID-19 Cases Using Past Viral Load
 Measurements. *Viruses* 14, (2022).

- 36. Sharmin, M. *et al.* Cross-sectional Ct distributions from qPCR tests can provide an early
 warning signal for the spread of COVID-19 in communities. *Front. Public Health* **11**, 1185720
 (2023).
- 37. Ahuja, V., Bowe, T., Warnock, G., Pitman, C. & Dwyer, D. E. Trends in SARS-CoV-2 cycle
- 659 threshold (Ct) values from nucleic acid testing predict the trajectory of COVID-19 waves.
- 660 *Pathology* **56**, 710–716 (2024).
- 38. Moro, A. *et al.* Trends in SARS-CoV-2 cycle threshold values in Bosnia and Herzegovina—A
 retrospective study. *Microorganisms* **12**, 1585 (2024).
- 39. Kissler, S. M. *et al.* Viral dynamics of SARS-CoV-2 variants in vaccinated and unvaccinated
 persons. *N. Engl. J. Med.* 385, 2489–2491 (2021).
- 40. Russell, T. W. *et al.* Combined analyses of within-host SARS-CoV-2 viral kinetics and
 information on past exposures to the virus in a human cohort identifies intrinsic differences
 of Omicron and Delta variants. *PLoS Biol.* 22, e3002463 (2024).
- 41. Fryer, H. R. *et al.* Viral burden is associated with age, vaccination, and viral variant in a
 population-representative study of SARS-CoV-2 that accounts for time-since-infectionrelated sampling bias. *PLoS Pathog.* **19**, e1011461 (2023).
- 42. Hay, J. A., Kennedy-Shaffer, L. & Mina, M. J. Viral loads observed under competing strain
 dynamics. *medRxiv* 2021.07.27.21261224 (2021).
- 43. Jones, T. C. *et al.* Estimating infectiousness throughout SARS-CoV-2 infection course
 Downloaded from. *Science* (2021) doi:10.1126/science.abi5273.
- 44. Wyllie, A. L. *et al.* Saliva or nasopharyngeal swab specimens for detection of SARS-CoV-2.
- 676 *N. Engl. J. Med.* **383**, 1283–1286 (2020).
- 45. Rhoads, D. *et al.* College of American pathologists (CAP) microbiology committee perspective: Caution must be used in interpreting the cycle threshold (ct) value. *Clin. Infect.*
- 679 *Dis.* **72**, e685–e686 (2021).

- 46. Arnaout, R. *et al.* The limit of detection matters: The case for benchmarking severe acute
 respiratory syndrome Coronavirus 2 testing. *Clin. Infect. Dis.* **73**, e3042–e3046 (2021).
- 47. Cuypers, L. et al. Nationwide harmonization effort for semi-quantitative reporting of SARS-
- 683 CoV-2 PCR test results in Belgium. *Viruses* **14**, 1294 (2022).
- 48. Hay, J. A. *et al.* Quantifying the impact of immune history and variant on SARS-CoV-2 viral
- kinetics and infection rebound: A retrospective cohort study. *Elife* **11**, (2022).
- 49. Elliott, P. *et al.* Rapid increase in Omicron infections in England during December 2021:
 REACT-1 study. *Science* 375, 1406–1411 (2022).
- 50. Brainard, J. *et al.* Comparison of surveillance systems for monitoring COVID-19 in England:
- a retrospective observational study. *Lancet Public Health* **8**, e850–e858 (2023).
- 51. Mellor, J. *et al.* Understanding the leading indicators of hospital admissions from COVID-19
- across successive waves in the UK. *Epidemiol. Infect.* **151**, e172 (2023).
- 692 52. https://www.mass.gov/info-details/covid-19-reporting.
- 53. http://dashboard.publichealth.lacounty.gov/covid19_surveillance_dashboard/.
- 54. Mgcv: Mixed GAM computation vehicle with automatic smoothness estimation.
 695 Comprehensive R Archive Network (CRAN) https://cran.r696 project.org/web/packages/mgcv/index.html.
- 55. Singanayagam, A. *et al.* Community transmission and viral load kinetics of the SARS-CoV-2
- delta (B.1.617.2) variant in vaccinated and unvaccinated individuals in the UK: a prospective,
- longitudinal, cohort study. *Lancet Infect. Dis.* **22**, 183–195 (2022).
- 56. lazymcmc. Preprint at https://github.com/jameshay218/lazymcmc.