

Sample Size Calculations for Variant Surveillance in the Presence of Biological and Systematic Biases

Shirlee Wohl^{1,2}, Elizabeth C. Lee¹, Bethany L. DiPrete^{3,4}, and Justin Lessler^{1,3,5}

¹ Johns Hopkins Bloomberg School of Public Health, Department of Epidemiology, Baltimore, MD, USA

² The Scripps Research Institute, Department of Immunology and Microbiology, La Jolla, CA, USA

³ Gillings School of Global Public Health, University of North Carolina at Chapel Hill, Department of Epidemiology, Chapel Hill, NC, USA

⁴ Injury Prevention Research Center, University of North Carolina at Chapel Hill, Chapel Hill, NC, USA

⁵ The Carolina Population Center, University of North Carolina at Chapel Hill, Chapel Hill, NC, USA

Correspondence to: swohl@scripps.edu

ABSTRACT

As demonstrated during the SARS-CoV-2 pandemic, detecting and tracking the emergence and spread of pathogen variants is an important component of monitoring infectious disease outbreaks. Pathogen genome sequencing has emerged as the primary tool for variant characterization, so it is important to consider the number of sequences needed when designing surveillance programs or studies, both to ensure accurate conclusions and to optimize use of limited resources. However, current approaches to calculating sample size for variant monitoring often do not account for the biological and logistical processes that can bias which infections are detected and which samples are ultimately selected for sequencing. In this manuscript, we introduce a framework that models the full process from infection detection to variant characterization and demonstrate how to use this framework to calculate appropriate sample sizes for sequencing-based surveillance studies. We consider both cross-sectional and continuous sampling, and we have implemented our method in a publicly available tool that allows users to estimate necessary sample sizes given a specific aim (e.g., variant detection or measuring variant prevalence) and sampling method. Our framework is designed to be easy to use, while also flexible enough to be adapted to other pathogens and surveillance scenarios.

MAIN TEXT

The emergence of SARS-CoV-2 variants with different epidemiologic properties has contributed to difficulties in controlling the COVID-19 pandemic. Towards the end of 2020, the first Variant of Concern (VOC) [1], later designated Alpha, was identified in the United Kingdom [2] and additional VOCs continued to be identified throughout the pandemic [3–5]. By definition, these variants are associated with increased COVID-19 transmissibility, increased virulence, or decreased effectiveness of available diagnostics, vaccines, or therapeutics, and many VOCs have had a significant impact on COVID-19 epidemiology [6]. Identifying these variants and tracking their spread in different populations is essential for informing our understanding of transmission and disease severity in this pandemic.

Whole genome sequencing of SARS-CoV-2 samples allows for detection of novel variants and regular monitoring of their frequency in populations. Although genomic sequencing has become faster and more cost-efficient, it is not possible or necessary to sequence all samples, and efficient allocation of resources (e.g., time and supplies) is critical in emergency public health responses. Therefore, selecting an appropriate sample for sequencing should play an important role in monitoring SARS-CoV-2 variants. However, sample sizes are still often dictated by cost and convenience, and there is limited guidance available for designing sampling strategies for genomic studies [7], including surveillance efforts aimed at detecting and characterizing VOCs.

Initial attempts at sample size calculations for tracking of SARS-CoV-2 VOCs have focused on determining the number of samples needed to identify a variant at a particular frequency in the population [8–11]. These frameworks often start at the point where COVID-19 samples are returned positive or start with assumptions about the proportion of COVID-19 infections that are tested and detected. However, the composition of this pool of detected infections may be affected by variant-specific differences in transmissibility, case detection, and test sensitivity. Hence, there is a need for a more comprehensive framework that models the full process underlying the selection of sequences. Furthermore, since samples are often sequenced regularly during an ongoing outbreak, sample size calculations should take into account changes in variant frequency over time.

Here we aim to develop an easy-to-use, actionable framework for selecting the appropriate number of SARS-CoV-2 samples for sequencing when the goal is (1) detecting new variants as they arise in a particular population or geographic region, or (2) measuring the prevalence of specific variants in that population. While our work focuses on genomic sequencing as the method for variant identification, this general framework can be applied to other molecular detection methods that differentiate between pathogen lineages.

Conceptual Framework and Approach

When designing a sequencing-based study or surveillance system, the first step is to identify the overall goal of the study, as different aims require different sample sizes in order to obtain reliable results. Whether variant detection or measuring prevalence is the primary goal [12], an important consideration is if samples are going to be collected in a single cross-sectional snapshot or through ongoing surveillance (**Fig 1**). This decision will influence both how the sample size is defined (i.e., overall study size versus average daily or weekly sampling rate), as well as what targets must be specified to calculate the appropriate sample size. For instance, in a cross-sectional study a possible target could be the probability of detecting a variant at a particular prevalence, while with ongoing surveillance it might be the waiting time to detect a recently introduced variant that is growing in prevalence. The requirements for specific goals are enumerated in **Figure 1**. If the study design is fixed (i.e., a set number of samples have already been collected or sequenced), the same principles can be applied to evaluate the questions that can be answered and confidence in the results.

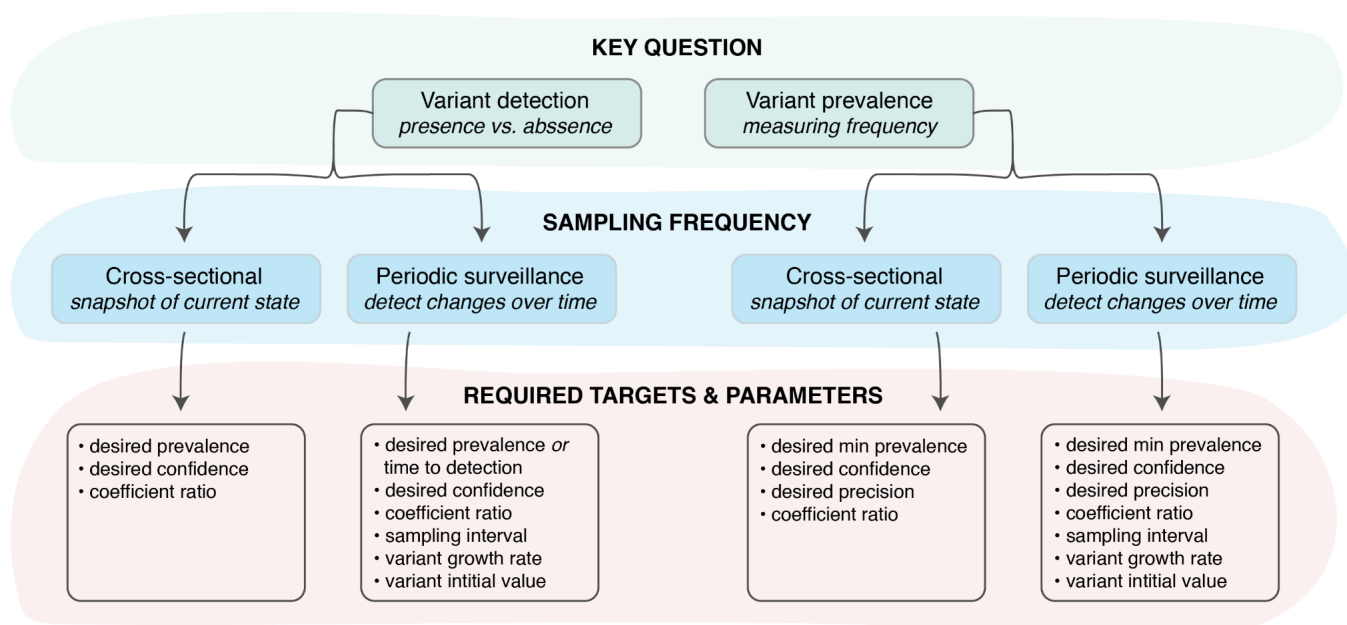


Figure 1. Decision tree for designing a variant surveillance program. Identifying the key goals of the study (green shaded region; variant detection or measuring variant prevalence) and sampling method (blue shaded region; cross-sectional or periodic) are necessary to determine the required targets and parameters (red shaded region) that must be specified in order to calculate the appropriate sample size.

Limitations of existing sampling strategies

For detection- and prevalence-based questions, we can use existing sampling theory as the basis of our approach. Specifically, the sample size needed to detect novel variants at some probability can be calculated with a simple application of the binomial distribution. The probability of detecting at least one case belonging to a variant of

interest (V_i) given the prevalence of this variant in the population (P_{V_i}) is equivalent to one minus the probability of not detecting it at all. Therefore, the sample size (n) needed to detect at least one case of a VOC at a pre-determined probability (p) has been shown to be [10]:

$$n = \frac{\log(1 - p)}{\log(1 - P_{V_i})} \quad (1)$$

Similarly, existing sample size theory can be used to estimate the prevalence of known VOCs in a population. Specifically, sample size calculations for estimating proportions can be used to determine the number of sequences that should be generated to estimate VOC prevalence within a desired confidence interval [13]:

$$n = \frac{Z^2 P_{V_i} (1 - P_{V_i})}{d^2} \quad (2)$$

Where n is the number of samples needed, Z is the Z-statistic for the 95% confidence level, P_{V_i} is the expected prevalence of the VOC in the population and d is the desired absolute precision (tolerance for error in the prevalence estimate). This methodology has previously been used to calculate the number of SARS-CoV-2 samples needed to detect variants at different frequency levels [8] and it assumes that the sample size (n) is small compared to the total infected population.

These approaches to sample size calculations are subject to limitations. For one, both equations assume that the pool of samples available for sequencing are a representative random sample of the total infected population. However, the biology and epidemiology of SARS-CoV-2 VOCs, such as severity of disease, may affect which samples are collected and sequenced (**Fig 2, Fig S1**). The sequences used for analysis may therefore not be directly reflective of the underlying distribution of viral sequences, and this bias may be detrimental or useful depending on the goals of surveillance.

Here, we characterize the mechanistic process from infection to case detection to variant identification using a simple modeling framework that captures how these processes differ between variants. We first explore how VOC attributes could bias detection of SARS-CoV-2 cases, and then determine how this bias may affect sample size calculations. We then extend the framework and focus on genomic surveillance as a continuous process, with sampling occurring periodically over time. Our framework has been implemented in a customizable, publicly available spreadsheet (<https://github.com/HopkinsIDD/VOCsamplesize>).

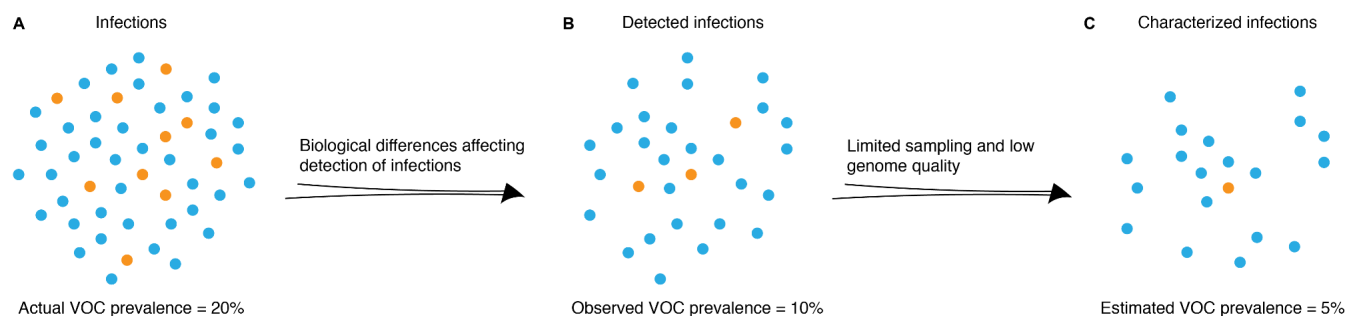


Figure 2. Factors affecting observed variant prevalence. VOC prevalence in (A) total population, (B) pool of detected infections, and (C) characterized infections (identified as a particular variant by sequencing or another technology). Biological and logistical differences between variants can lead to bias in observed variant proportions. Orange = infections caused by VOC (variant of concern); blue = infections caused by other variants of the same pathogen.

Variant Surveillance Model

As discussed above, we aimed to characterize the factors that could affect the collection of SARS-CoV-2 samples and their selection for downstream processes such as sequencing. To do this, we developed a model that tracks how biological differences between variants, as well as logistical challenges in case and variant detection, may affect estimated variant frequency, given a true underlying frequency in a population. This model distinguishes all infections (N), from those detected positive (D), from those that produce a genome sequence (G) that can be used to identify the underlying variant.

We conceive of the model in two phases: 1) *infection detection*, which describes the joint biological and testing mechanisms that lead infections (N) to be detected as COVID-positive (D) by a surveillance system (Fig 3, top row), and 2) *infection characterization*, which describes the selection of samples for genomic sequencing from high quality detected infections (H) and identification of specific variants from the resulting high-quality sequences (Fig 3, bottom row). In this context, “high-quality detected infections” refers to pathogen-positive samples of high enough quality (e.g., by a metric such as cycle threshold value) that they will be selected for sequencing, while “high-quality sequences” refers to pathogen sequences that are complete enough to characterize the infection-causing variant.

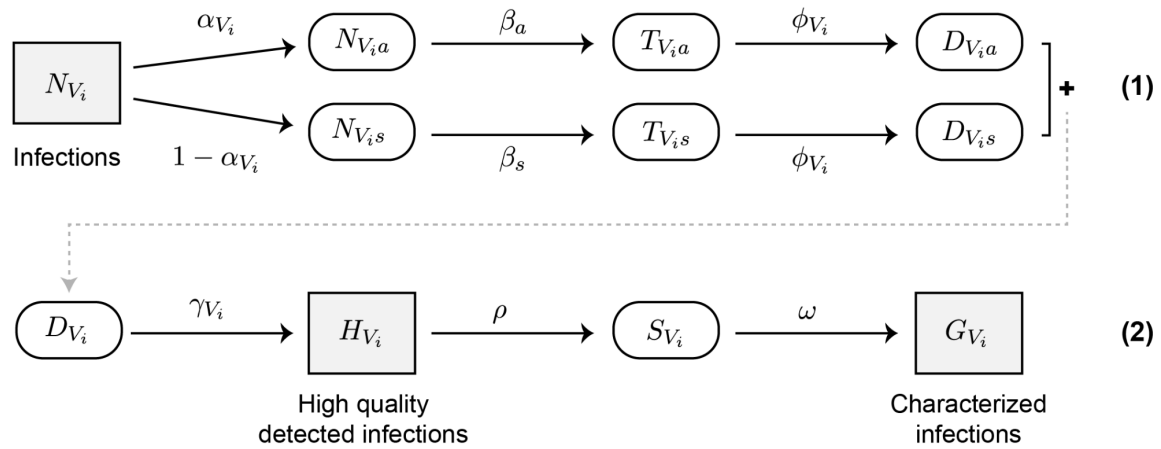


Figure 3. Schematic of variant surveillance model. (1) infection detection process; (2) infection characterization process. Parameters are defined in **Table 1**.

Model states are separated by transition parameters that model how biological differences between variants can affect factors such as testing rates, testing sensitivity, and sample quality (**Table 1**). The model is generalizable for any number of variants of interest, where the variant(s) of interest are always compared to the remaining population. For example, the number of infections caused by a single variant of interest V_1 is tracked alongside the number of infections not caused by this variant, which we term V_2 . This setup allows us to track the proportion of variant i at any given step, which is often more interesting than the total number.

Using this model, the number of high quality detected infections attributable to a specific variant is as follows:

$$H_{V_i} = NP_{V_i} \phi_{V_i} \gamma_{V_i} (\alpha_{V_i} \beta_a + (1 - \alpha_{V_i}) \beta_s) \quad (3)$$

By calculating this quantity for each variant of interest and the remainder of the population, we can determine the prevalence of each VOC in the pool of high quality samples available for sequencing. When the detection parameters do not vary between variants, the VOC prevalence in detected high quality infections mirrors its prevalence in the greater population. Variation in these parameters between pathogen variants can be summarized in a single parameter, which we term the *coefficient of detection*:

$$C_{V_i} = \phi_{V_i} \gamma_{V_i} (\alpha_{V_i} \beta_a + (1 - \alpha_{V_i}) \beta_s) \quad (4)$$

The value of this coefficient for each variant will determine the bias already present in H , the population from which we ultimately draw our sample.

Table 1. Model states and parameters.

Param	Description	Param	Description
N	Total number of infections in population	H_{V_i}	Number of detected, high quality, infections caused by variant i
P_{V_i}	Proportion of variant i in population	S_{V_i}	Number of detected, high quality, infections caused by variant i that are selected for sequencing
N_{V_i}	Number of infections caused by variant i	G_{V_i}	Number of high quality sequences from infections caused by variant i
$N_{V_i^a}$	Number of asymptomatic infections caused by variant i	α_{V_i}	Probability that an infection caused by variant i is asymptomatic
$N_{V_i^s}$	Number of symptomatic infections caused by variant i	β_x	Probability of testing, given type of infection (x : symptomatic or asymptomatic)
$T_{V_i^a}$	Number of tested asymptomatic infections caused by variant i	ϕ_{V_i}	Probability that a tested infection caused by variant i results in a positive test (sensitivity)
$T_{V_i^s}$	Number of tested symptomatic infections caused by variant i	γ_{V_i}	Probability that a detected infection caused by variant i meets some quality threshold
$D_{V_i^a}$	Number of detected asymptomatic infections caused by variant i	ρ	Probability that a sample is selected for sequencing
$D_{V_i^s}$	Number of detected symptomatic infections caused by variant i	ω	Probability that a sequenced sample produces a high quality genome
D_{V_i}	Number of detected infections (symptomatic and asymptomatic) caused by variant i		

Effects of pathogen properties on variant surveillance

We explored the effects of biological and logistical differences between variants—summarized in the coefficient of detection (**Equation 4**)—on the variant proportions observed in H , the pool of high quality detected infections from which to sample (**Equation 3**). We can calculate the multiplicative bias in our estimate of a particular variant as a function of the underlying prevalences and coefficients of detection for all variants in the system:

$$\frac{\text{actual } V_1 \text{ prevalence}}{\text{observed } V_1 \text{ prevalence}} = P_{V_1} + \frac{C_{V_2}}{C_{V_1}} P_{V_2} + \frac{C_{V_3}}{C_{V_1}} P_{V_3} + \dots + \frac{C_{V_n}}{C_{V_1}} P_{V_n} \quad (5)$$

Where n is the total number of variants in the population ($n \geq 2$).

Unsurprisingly, a larger differential between the coefficient of detection for V_1 and the coefficients of detection for other variants in the system leads to more bias in the observed frequency. Additionally, the observed prevalence of V_1 in H is more biased when P_{V_1} is smaller (**Fig 4A**).

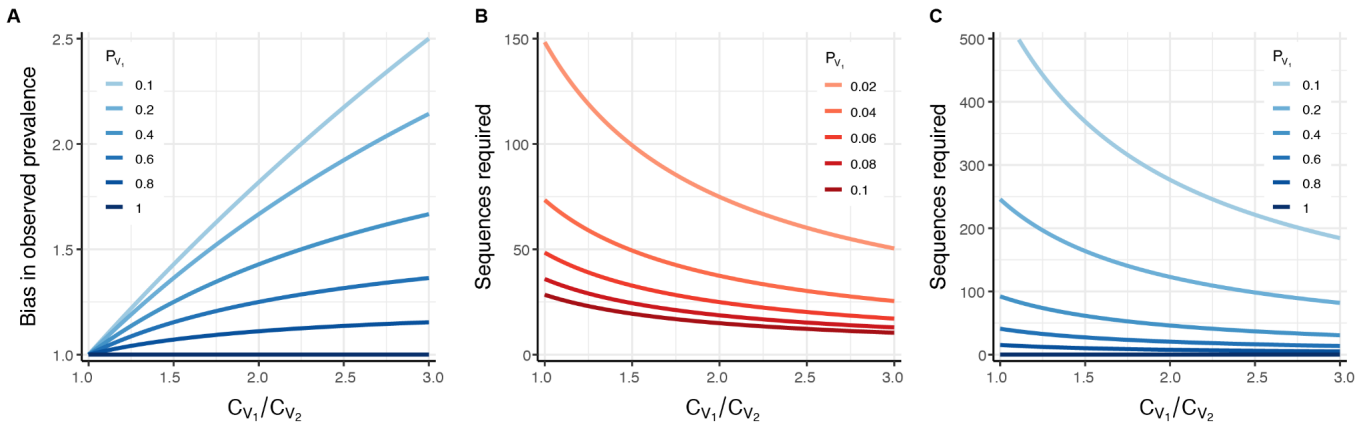


Figure 4. Exploring the effects of the coefficient of detection. (A) Multiplicative bias in the observed prevalence of variant V_1 in H , the pool of high quality infected detections to sample from (bias: observed V_1 prevalence divided by actual V_1 prevalence). (B) Number of sequences required to detect at least one infection caused by V_1 with 95% probability, for different V_1 prevalence values and coefficient of detection ratios. (C) Number of sequences required to determine the prevalence of variants with a frequency of at least P_{V_1} in the population, with 95% confidence and 25% precision. The prevalence calculated with these sequences will reflect the observed (biased) value, and will need to be corrected using **Equation 6**. All panels assume a two-variant system, where V_1 is the variant of interest and V_2 is the rest of the pathogen population. In (B) and (C), note that the number of samples selected for sequencing should exceed the number of sequences required if $\omega < 1$.

We can also calculate a correction factor q such that:

$$P_{V_1} = \frac{q(\text{odds}_{V_1}^*)}{1 + q(\text{odds}_{V_1}^*)} \quad (6)$$

Where $\text{odds}_{V_1}^*$ is the observed odds of the V_1 prevalence in H . This equation allow for a direct conversion between the observed variant frequencies in the sampling pool ($P_{V_i}^*$) and the true frequency of V_1 in the infected population (P_{V_1}). In a two-variant system (i.e., a system with one variant of interest compared to the rest of the population),

$$q = \frac{C_{V_2}}{C_{V_1}} \quad (\text{see **Appendix** for derivation and correction factor values in larger systems}).$$

Sampling strategies for cross-sectional surveillance

In the following sections, we provide examples of how to calculate the appropriate sample size for surveillance given potential biases in observed variant frequencies. We also discuss how this bias—or, in some cases, enrichment—may make it easier to detect or measure the prevalence of certain variants.

Variant detection

Detecting the introduction of new variants into specific populations is a common goal during a pathogen outbreak. This requires identification of variants while they are still at low frequency in the population. For example, we may be interested in determining the minimum sample size needed to have a 95% chance of detecting a variant at 2% frequency in a specific population. If this variant is biologically and epidemiologically identical to the rest of the population, its frequency in the sampling pool (H) will reflect its frequency in the overall population. In this case, we can apply binomial sampling theory (**Equation 1**) to calculate the number of sequences needed (see **Fig S2** for validation of binomial sampling process):

$$n = \frac{\log(1 - p)}{\log(1 - P_{V_1})} = \frac{\log(1 - 0.95)}{\log(1 - 0.02)} = 149$$

Most variants of interest, however, are not biologically and epidemiologically identical to the rest of the pathogen population. For example, variants can emerge that are more transmissible, such as the SARS-CoV-2 Delta variant [1,14,15]. This increased transmissibility can be for a variety of reasons, such as higher viral loads in infected patients or more efficient entry into host cells, all of which require adjustment to these calculations. Here, we assume that a VOC is more transmissible specifically because it causes higher viral titers in infected patients, and that this increased titer increases the testing sensitivity of the variant ($\phi_{V_1} = 0.975$) as compared to the rest of the population ($\phi_{V_2} = 0.95$). We also assume that detected infections caused by this VOC contain more virus and therefore have an increased probability ($\gamma_{V_1} = 0.8$) of meeting quality thresholds (e.g., Ct value cutoffs) than other positive samples ($\gamma_{V_2} = 0.6$). We assume that all other biological and surveillance parameters are the same between the VOC and the rest of the pathogen population.

Using these parameters, we can calculate the ratio of coefficients of detection and use this to calculate the VOC frequency we expect to see in our sample. Rearranging **Equation 5**, we see that:

$$\text{observed } V_1 \text{ prevalence} = \frac{P_{V_1}}{P_{V_1} + \frac{C_{V_2}}{C_{V_1}} P_{V_2}} = \frac{0.02}{0.02 + \left(\frac{0.95 \cdot 0.6}{0.975 \cdot 0.8}\right)(0.98)} = 0.027$$

We then apply sampling theory as above, using the observed variant frequency (2.7%) as P_{V_i} . Because the variant is enriched in our population of detected infections, we find that only 110 sequences are needed to be 95% confident in detection of this variant (**Fig 4B**; see **Fig S3A** for sequence requirements for 50% confidence). Since not every sequenced sample produces a usable sequence, even after selecting for high quality samples (**Fig 3**), we assume a sequencing success rate of 80% for all variants ($\omega = 0.8$), which means 138 samples should be selected for sequencing in order to obtain 110 complete genomes. The same procedure could be performed to determine the

sample size needed to detect a more severe variant—or any variant that has some effect on pathogen detection—provided the ratio of coefficients of detection can be estimated.

Variant prevalence

After a variant is first detected, sequencing is often used to monitor its frequency in the population, and to note any increases in variant frequency that may suggest epidemiological or biological trends. Therefore, we assume that we are interested in calculating the minimum sample size needed to correctly (within 25% of the true value) determine the prevalence of a variant at >10% frequency in the population with 95% confidence. If this variant is biologically and epidemiologically identical to the rest of the population, its frequency in the sampling pool (H) will reflect its frequency in the population. In this case, we can apply existing theory (**Equation 2**) to calculate the number of sequences needed:

$$n = \frac{Z^2 P_{V_1} (1 - P_{V_1})}{d^2} = \frac{1.96^2 (0.1)(1 - 0.1)}{(0.1 * 0.25)^2} = 554$$

In this example, we calculate the sample size with the smallest prevalence (10%) we are interested in accurately measuring, since this requires the largest sample size. We do not apply any sort of finite population size correction [8], though this could decrease the sample size needed for prevalence estimation.

If a variant of interest has differing biological or epidemiological properties, we must adjust our calculations to account for the likely over or underrepresentation of this variant in the sampling pool. A more severe variant, for example, may decrease the proportion of infected individuals who are asymptomatic ($\alpha_{V_1} = 0.25$), compared to the rest of the population ($\alpha_{V_2} = 0.4$), but may have a limited effect on the other biological and surveillance parameters. A difference in the asymptomatic rate only biases the observed variant frequency if testing rates are different for symptomatic and asymptomatic infections (see **Equation 4**), so in this example we assume that symptomatic infections are tested at a higher frequency ($\beta_s = 0.3$) than asymptomatic infections ($\beta_a = 0.05$). We again use **Equations 4 and 5** to calculate the variant frequency we expect in our sample, assuming a true underlying frequency of 10%:

$$\text{observed } V_1 \text{ prevalence} = \frac{P_{V_1}}{P_{V_1} + \frac{c_{V_2}}{c_{V_1}} P_{V_2}} = \frac{0.01}{0.01 + (0.84)(0.9)} = 0.117$$

We then apply sampling theory as above, using $P_{V_i} = 0.117$. We find that the enrichment of VOC samples among detected infections means fewer sequences ($n = 464$) are needed for measuring the prevalence of this variant (**Fig 4C**; see **Fig S3B-C** for sequence requirements with different confidence and precision values). If we assume a sequencing success rate of 80%, a total of 580 samples should be selected for sequencing in this example. If

resources are constrained or not enough samples can be collected, the confidence level for both monitoring and detection can be calculated from the number of samples available.

When later using these sequences to estimate VOC prevalence, it is important to note that the variant prevalence estimated from the sequence (or other variant characterization) data will be the *observed variant prevalence*, even if the required number of sequences are available. In other words, because the sampling pool itself is biased, the proportion of sequences that are characterized as V_1 is equal to $P_{V_1}^*$ (and not the true underlying prevalence, P_{V_1}), even if every available sample is sequenced. **Equation 6** must then be applied to estimate the true variant prevalence from this observed value.

Sampling strategies for ongoing surveillance

During an infectious disease outbreak, variant detection and monitoring are ongoing processes. Focusing on sampling strategies over time will allow us to answer more realistic questions, including: what sample size is required to ensure detection of a variant in a pre-specified amount of time, or before the variant reaches a particular prevalence in the population?

Variant detection

To answer these questions, we assume that the same number of sequences are sampled at each time step. Even if samples are not sequenced at each time step (e.g., every day), we assume that the same number of samples from each day are included in a (e.g., weekly) sequencing batch, effectively increasing sequencing frequency to daily with some delay in final results. Given these assumptions, we can again use binomial sampling theory (**Equation 1**) to calculate the probability of detecting a VOC on or before time step t . The resulting equation takes the form of a survival function, as follows:

$$\Pr(d \leq t) = 1 - \prod_{x=0}^{x=t} (1 - P_x)^n$$

Where $\Pr(d \leq t)$ is the probability of detection on or before time t , n is the sample size per unit time, and P_x is the prevalence of the variant of interest in the population at time x . After rearranging this equation to solve for the per-timestep sample size and approximating the product with a continuous function (see **Appendix** for full derivation), we obtain:

$$n = -\frac{\ln(1 - \Pr(d \leq t))}{G(t) - G(0)} \quad (7)$$

Where $G(t)$ is the cumulative density of the function used to model variant growth over time, at time t . In other words, we can easily estimate the necessary sample size per time unit, provided we can approximate how the variant prevalence is changing over time.

In the example below, we assume that variant prevalence follows a logistic growth curve. Logistic growth is often ascribed to variants with a fitness advantage, such as the Alpha SARS-CoV-2 variant [16], though in this section we will assume that the variant of interest does not affect any of the parameters that go into calculating the coefficient of detection (we will relax this assumption in the following section). We assume there was a single introduction of this variant into a population of 10,000 individuals and that the growth rate is approximately 0.1 per day [17]. Now, we can use **Equation 7** to calculate the per-day sample size needed to ensure detection (with 95% probability) of this Alpha-like variant within 14 days of its initial emergence:

$$n = -\frac{\ln(1 - 0.95)}{G(14) - G(0)} = 158 \text{ sequences/day}$$

In other words, generating $158 * 7 = 1,106$ sequences per week (assuming sequences are well-distributed throughout the week) ensures a 95% probability of detection of this variant within 14 days of initial introduction. It is important to note that, given a single introduction and a growth rate of 0.1 per day, the variant will only be at a frequency of 0.04% by day 14. It may be more realistic to assume multiple introductions of a highly transmissible variant, or to compute the sample size needed to detect the variant before it reaches a particular prevalence in the population. If we instead assume 3 introductions into a population of 10,000 and the same growth rate of 0.1 per day, we can calculate that the prevalence will surpass 1% on day 36. We can again apply **Equation 7** to calculate the number of sequences needed per day to ensure detection by the time the VOC prevalence surpasses 1%:

$$n = -\frac{\ln(1 - 0.95)}{G(36) - G(0)} = 29 \text{ sequences/day}$$

Generating 29 sequences per day, or just over 200 sequences per week, may be a much more manageable number. That said, it is important to consider the sequencing success rate (e.g., $\omega = 0.8$) when calculating the number of samples that should be selected for sequencing. To generate 203 high quality sequences per week, 254 samples will need to be selected for sequencing.

Variant detection with a biased sample

As discussed above, VOC prevalence may be enriched in the sampling pool, meaning that fewer sequences may be needed for confident detection of the variant. Using **Equation 5**, we can calculate the observed variant frequency at each time step given a growth rate and starting variant prevalence (e.g., 3 introductions into a population of 10,000) for a two-variant system as follows:

$$\text{observed } V_1 \text{ prevalence} = \frac{g(t)}{g(t) + \frac{C_{V_2}}{C_{V_1}}(1 - g(t))}$$

Where $\frac{C_{V_2}}{C_{V_1}}$ represents the relative coefficients of detection between the general pathogen population (C_{V_2}) and the variant of interest (C_{V_1}). Furthermore, $g(t)$ is the probability density of the function used to model variant growth over time, at time t . From this, we can calculate the cumulative density function of the *observed* variant prevalence distribution, and use this approximation in our sample size calculation.

If we assume that the Alpha-like VOC described above results in a coefficient of detection ratio of $\frac{C_{V_2}}{C_{V_1}} = \frac{0.95 \cdot 0.6}{0.975 \cdot 0.8}$ (see ‘Variant Detection’ in ‘Sampling Strategies for Cross-Sectional Surveillance’), the sample size needed to ensure 95% of detection by the time the VOC prevalence surpasses 1% in the population is:

$$n = -\frac{\ln(1 - 0.95)}{G^*(36) - G^*(0)} = 21 \text{ sequences/day}$$

Where $G^*(t)$ is the cumulative density of the function used to model the growth of the *observed* variant frequency over time, at time t , and we again assume an initial prevalence of 3 in 10,000 and a growth rate of 0.1 per day (see **Appendix** for full derivation). As expected, the enrichment of the VOC in the sampling pool decreases the number of sequences needed for detection (**Fig 5**). We can also use **Figure 5** (and **Fig S4**) to evaluate the marginal costs and benefits of changing the number of samples selected for sequencing, which may allow for the design of surveillance systems that take maximum advantage of available resources.

Implementation

We have implemented the method described above in a publicly available spreadsheet (<https://github.com/HopkinsIDD/VOCsamplesize>). This spreadsheet can be used to calculate the required sample size in each of the three scenarios described above: cross-sectional sampling for variant detection, cross-sectional sampling for measuring variant prevalence, and periodic sampling for variant detection. The equations implemented in this spreadsheet can be used both backwards and forwards—a user can input epidemiological and biological parameters and use them to determine the sample size needed to achieve the primary aim (detection or measuring prevalence) given a desired confidence level, or they can input a desired sample size and use them to calculate confidence in the results.

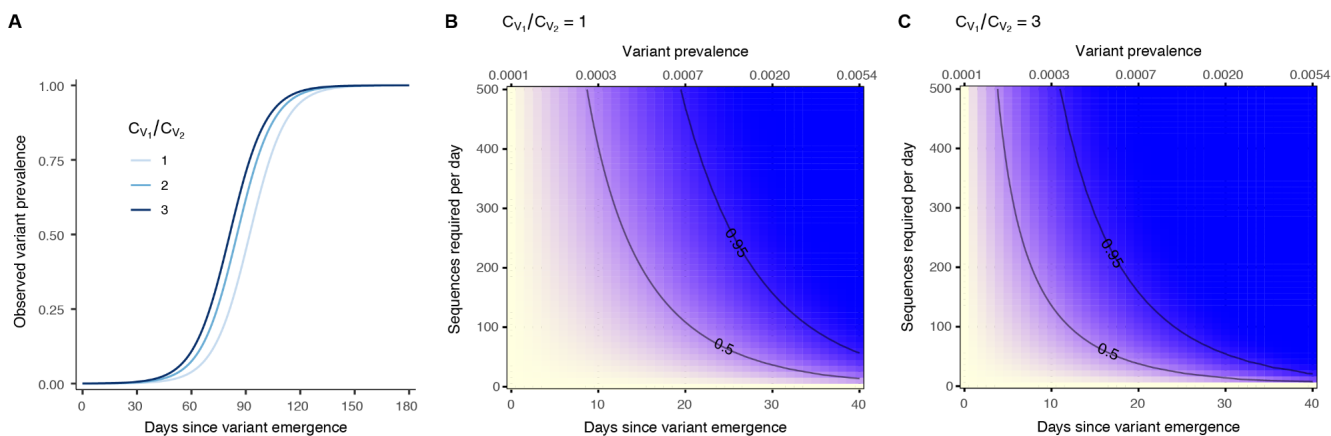


Figure 5. Sample size required for detection of a variant growing in prevalence. For a variant whose prevalence increases following a logistic curve with growth rate = 0.1 per day and starting value = 1/10000: **(A)** Observed variant prevalence over time given different coefficient-of-detection ratios. **(B-C)** Probability of detecting at least one infection caused by V_1 (yellow = 0% probability; blue = 100% probability) on or before a specific day (bottom x-axis) or desired prevalence (top x-axis), given per-day sample size and specified coefficient of detection ratio. Note that the desired prevalence (top x-axis) is the actual variant prevalence in the population and that the number of samples selected for sequencing should exceed the number of sequences required if $\omega < 1$. 50% and 95% probability of detection contours are indicated.

Discussion

Designing pathogen surveillance systems must begin with identifying the primary purpose or key questions to be answered with the system. For example, surveillance strategies will be different when the goal is early detection of a newly introduced VOC versus when the goal is measuring the prevalence of an existing variant [12]. In either case, there are a myriad of factors that influence which infections are ultimately sequenced. Here we present a framework for thinking about these factors, and we show that their effects can be summarized in a single number, the coefficient of detection. This coefficient characterizes how biological and logistical factors can lead to VOC enrichment in a sample—leading to earlier detection—while also biasing measurement of the true underlying VOC prevalence. Depending on the purpose of surveillance, it will be important to account for these effects in sample size calculations and subsequent reporting of results. The work presented here aims to provide an accessible set of methods for doing so, and a general approach that can be extended to other settings and study designs. In addition to providing evidence-based guidance for sampling design, our framework can be applied retrospectively to evaluate the accuracy of a detection or variant prevalence result based on the number of samples sequenced.

A perceived barrier to using the approach outlined here may be lack of knowledge of the exact parameters that are summarized by the coefficient of detection. However, it is not necessary to know individual parameter values when using our framework as long as we can approximate their ratio, since all calculations rely solely on the ratio of coefficients of detection. Likewise, parameters that have the same value across variants need not be specified at all.

Although decreasing the number of parameters that need to be specified makes the framework easier to use, there is still value in breaking down the process of surveillance into its component parts. For instance, it may be difficult to estimate a single pathogen testing rate in settings without consistent testing of asymptomatic individuals (e.g., hospitals or settings with limited testing capacity). But if the asymptomatic and symptomatic testing rates are separated into two parameters, we can assume the asymptomatic testing rate is negligible (or at least similar between variants) and focus on the symptomatic testing rate, which may be easier to quantify.

In considering the full process from infection detection to variant characterization, we have aimed to make our framework flexible enough to handle situations not explicitly discussed above. Although most of the examples presented in this manuscript focus on a two-variant model, the framework is set up to allow for exploration of multiple variants simultaneously (see **Appendix**). Further, while we focus on detecting variants undergoing logistic growth, the sample size needed to detect a variant can be calculated for any growth function as long as the functional form and underlying parameters are known. Additionally, while we have tried to identify the key processes that affect pathogen detection, the coefficient of detection could be modified to incorporate other parameters that differ between variants, or to include factors that may affect which sequences produce complete genomes (i.e., factors that affect the variant characterization process shown in the bottom part of **Fig 3**). For example, we could allow the sequencing success rate (ω) to differ between variants (e.g., due to differences in primer binding when using PCR-based methods for variant characterization or amplification prior to sequencing), despite the use of an initial sample quality filter (γ). Finally, the framework could be extended to any pathogen for which there is some method (e.g., variant-specific PCR assays) to differentiate pathogen lineages with potentially different epidemiological or biological processes.

While this manuscript covers sample size calculations for variant detection given both cross-sectional and periodic sampling approaches, additional work is needed to determine appropriate sample sizes for measuring VOC prevalence with periodic sampling. This would enable interpretation of small changes in variant prevalence over time that could serve as an early indication of changing epidemiological dynamics in the population. Additionally, when multiple variants are present in the population, accurate prevalence estimation of one VOC necessarily constrains the potential prevalence values of another VOC; expanding the framework to co-estimate prevalences for multiple VOCs may make it possible to leverage this interdependence and further reduce the sample sizes required for accurate monitoring.

When designing surveillance systems based on this framework, it is important to remember that infection and sampling processes can be heterogeneous in ways not captured by our model. Future work could consider the effects of spatial heterogeneity on transmission and sampling, the impact of time-varying model parameters in the continuous surveillance context, or could extend the framework to metapopulations. Given the current framework and assumption of homogeneity, samples selected for sequencing should be selected as randomly as possible, or selected in a way that maximizes the geographic and temporal distribution of sequences.

Whole genome sequencing has revealed the importance of characterizing and monitoring specific pathogen variants during an ongoing epidemic. As this technology becomes more accessible and central to our understanding of established and emerging pathogens, it is important that we improve the rigor with which we design studies that use this data. Sophisticated modeling approaches have been invaluable in improving how we collect and interpret pathogen genomic information, but most are neither nimble nor accessible enough to be widely used during a crisis. Similarly, ad-hoc approaches or classical study designs may not lead to the optimal allocation of resources. Here we have attempted to lay out a framework that is widely accessible yet still accounts for many of the factors that uniquely impact pathogen genomic studies and surveillance programs. As the SARS-CoV-2 pandemic continues and new infectious threats arise, we hope this approach will help better guide the collection of data that has proven critical to the pandemic response and serve as a starting point for further methodological innovation.

ACKNOWLEDGMENTS

We thank Edyth Parker for her insightful comments on the manuscript. Funding was provided by Bill and Melinda Gates Foundation INV-025321 (S.W.) and OPP1195157 (S.W. and J.L.).

REFERENCES

1. Tracking SARS-CoV-2 variants. [cited 23 Aug 2021]. Available: <https://www.who.int/en/activities/tracking-SARS-CoV-2-variants/>
2. Public Health England. Investigation of novel SARS-CoV-2 variant: Variant of Concern 202012/01. Available: https://assets.publishing.service.gov.uk/government/uploads/system/uploads/attachment_data/file/959438/Technical_Briefing_VOC_SH_NJL2_SH2.pdf
3. Tegally H, Wilkinson E, Giovanetti M, Iranzadeh A, Fonseca V, Giandhari J, et al. Emergence and rapid spread of a new severe acute respiratory syndrome-related coronavirus 2 (SARS-CoV-2) lineage with multiple spike mutations in South Africa. *bioRxiv. medRxiv*; 2020. doi:10.1101/2020.12.21.20248640
4. Faria NR, Mellan TA, Whittaker C, Claro IM, Candido D da S, Mishra S, et al. Genomics and epidemiology of the P.1 SARS-CoV-2 lineage in Manaus, Brazil. *Science*. 2021. doi:10.1126/science.abh2644
5. Classification of Omicron (B.1.1.529): SARS-CoV-2 Variant of Concern. [cited 21 Dec 2021]. Available: [https://www.who.int/news/item/26-11-2021-classification-of-omicron-\(b.1.1.529\)-sars-cov-2-variant-of-concern](https://www.who.int/news/item/26-11-2021-classification-of-omicron-(b.1.1.529)-sars-cov-2-variant-of-concern)
6. Bushman M, Kahn R, Taylor BP, Lipsitch M, Hanage WP. Population impact of SARS-CoV-2 variants with enhanced transmissibility and/or partial immune escape. *Cell*. 2021. doi:10.1016/j.cell.2021.11.026
7. Wohl S, Giles JR, Lessler J. Sample size calculation for phylogenetic case linkage. *PLoS Comput Biol*. 2021;17:e1009182.
8. European Centre for Disease Prevention and Control. Sequencing of SARS-CoV-2: first update. 18-January-2021. Available: <https://www.ecdc.europa.eu/sites/default/files/documents/Sequencing-of-SARS-CoV-2-first-update.pdf>
9. Vavrek D, Speroni L, Curnow KJ, Oberholzer M, Moeder V, Febbo PG. Genomic surveillance at scale is required to detect newly emerging strains at an early timepoint. *bioRxiv. medRxiv*; 2021. doi:10.1101/2021.01.12.21249613
10. The University of Texas COVID-19 Modeling Consortium. Sample Size Calculator Detecting COVID-19 Variants. In: Variant Detection Calculator [Internet]. [cited 1 May 2021]. Available: <https://covid-19.tacc.utexas.edu/dashboards/variants/>
11. Brito AF, Semenova E, Dudas G, Hassler GW, Kalinich CC, Kraemer MUG, et al. Global disparities in SARS-CoV-2 genomic surveillance. *medRxiv*. 2021. doi:10.1101/2021.08.21.21262393
12. European Centre for Disease Prevention and Control. Guidance for representative and targeted genomic SARS-CoV-2 monitoring. 2021 May. Available: <https://www.ecdc.europa.eu/sites/default/files/documents/Guidance-for-representative-and-targeted-genomic-SARS-CoV-2-monitoring.pdf>
13. Wayne W. Daniel CLC. *Biostatistics: A Foundation for Analysis in the Health Sciences*, 11th Edition. Wiley; 2018.
14. Liu Y, Rocklöv J. The reproductive number of the Delta variant of SARS-CoV-2 is far higher compared to the ancestral SARS-CoV-2 virus. *J Travel Med*. 2021;28. doi:10.1093/jtm/taab124
15. Challen R, Dyson L, Overton CE, Guzman-Rincon LM, Hill EM, Stage HB, et al. Early epidemiological signatures of novel SARS-CoV-2 variants: establishment of B.1.617.2 in England. *bioRxiv. medRxiv*; 2021. doi:10.1101/2021.06.05.21258365
16. Volz E, Mishra S, Chand M, Barrett JC, Johnson R, Geidelberg L, et al. Assessing transmissibility of SARS-CoV-2 lineage B.1.1.7 in England. *Nature*. 2021;593: 266–269.
17. Davies NG, Abbott S, Barnard RC, Jarvis CI, Kucharski AJ, Munday JD, et al. Estimated transmissibility and impact of SARS-CoV-2 lineage B.1.1.7 in England. *Science*. 2021;372. doi:10.1126/science.abg3055

SUPPLEMENTAL FIGURES

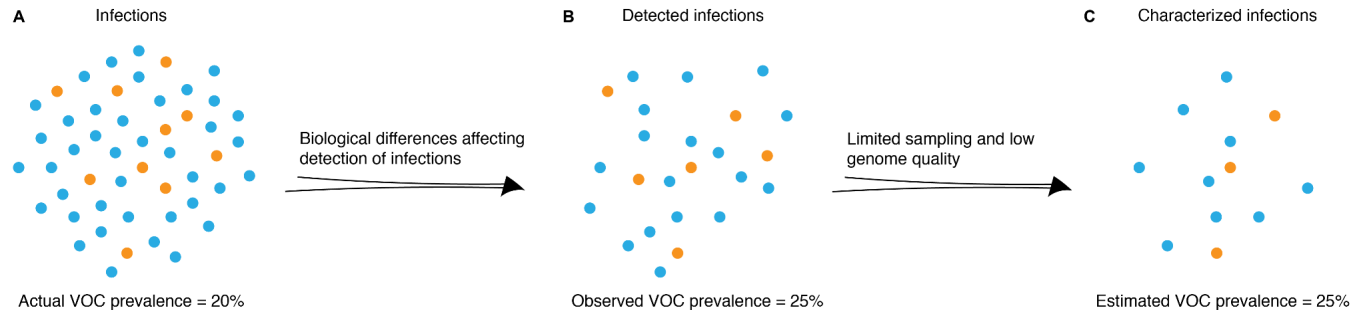


Figure S1. Factors resulting in enrichment of observed variant prevalence. VOC prevalence in (A) total population, (B) pool of detected infections, and (C) characterized infections (identified as a particular variant by sequencing or another technology). Biological between variants can lead to enrichment of VOC in observed variant proportions. Orange = infections caused by VOC (variant of concern); blue = infections caused by other variants of the same pathogen.

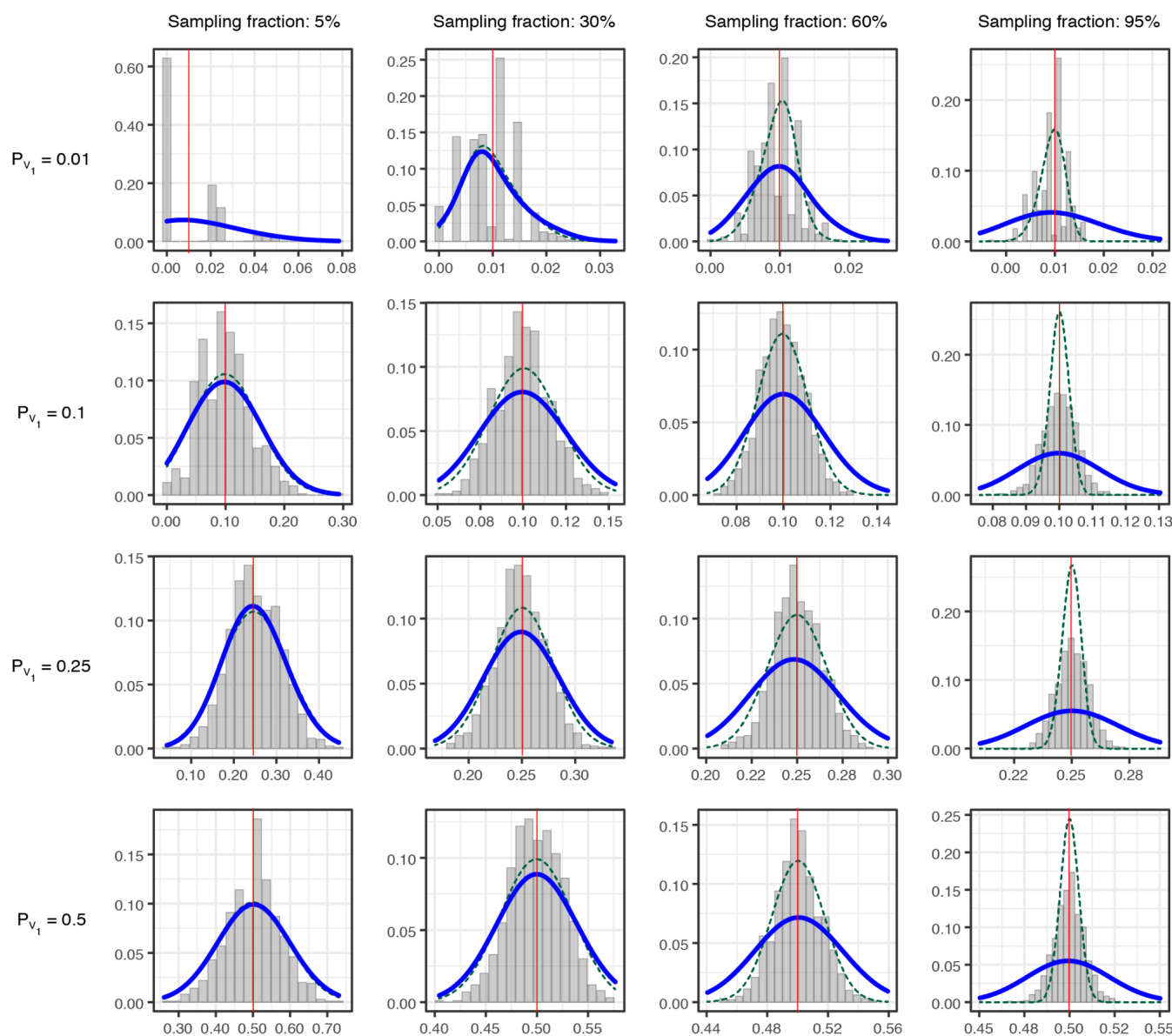


Figure S2. Validation of binomial sampling approximation. Variant prevalence estimates (x-axis) from 1,000 simulations of our model for each combination of variant prevalence (P_{V_1}) (rows) and sampling proportion (columns). Simulations each assume 10,000 infected individuals in the population and $C_{V_1} = C_{V_2} = 0.114$, resulting in 1,140 high quality detected infections (H) from which to sample. In each simulation, $\omega = 0.8$ and infections that progress between model states are selected stochastically using a binomial process. Blue lines = binomial distribution (i.e., sampling with replacement, an approximation of the simulated sampling process) given stated sampling fraction (from H) and variant prevalence; green dotted lines = hypergeometric distribution (i.e., sampling without replacement, the exact sampling process) given stated sampling fraction and variant prevalence; red vertical lines = marker of input variant frequency. The binomial distribution approximates the simulated process well, except when nearly all detected infections are sampled (as expected) or the variant prevalence is very low.

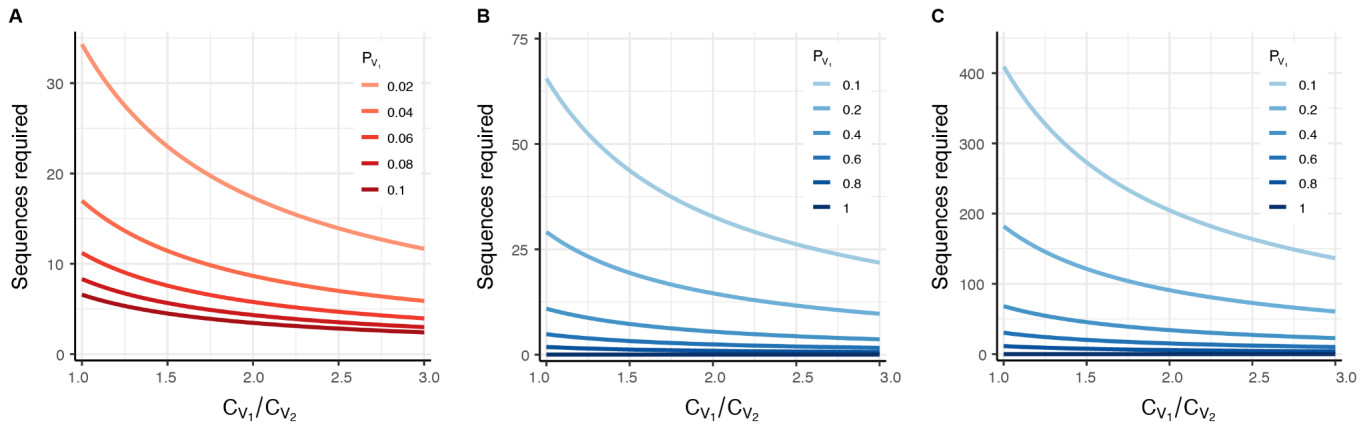


Figure S3. Sample size needed for variant detection and prevalence estimation with 50% confidence. (A) Number of sequences required to detect at least one infection caused by V_1 with 50% probability, for different V_1 prevalence values and coefficient of detection ratios. **(B)** Number of sequences required to determine the prevalence of variants with a frequency of at least P_{V_1} in the population, with 50% confidence and 25% precision. **(C)** Same as (B), but with 50% confidence and 10% precision. The calculated prevalence in (B) and (C) will reflect the observed (biased) value, and will need to be corrected using **Equation 6**. These figures assume a two-variant system, where V_1 is the variant of interest and V_2 is the rest of the pathogen population. In all panels, note that the number of samples selected for sequencing should exceed the number of sequences required if $\omega < 1$.

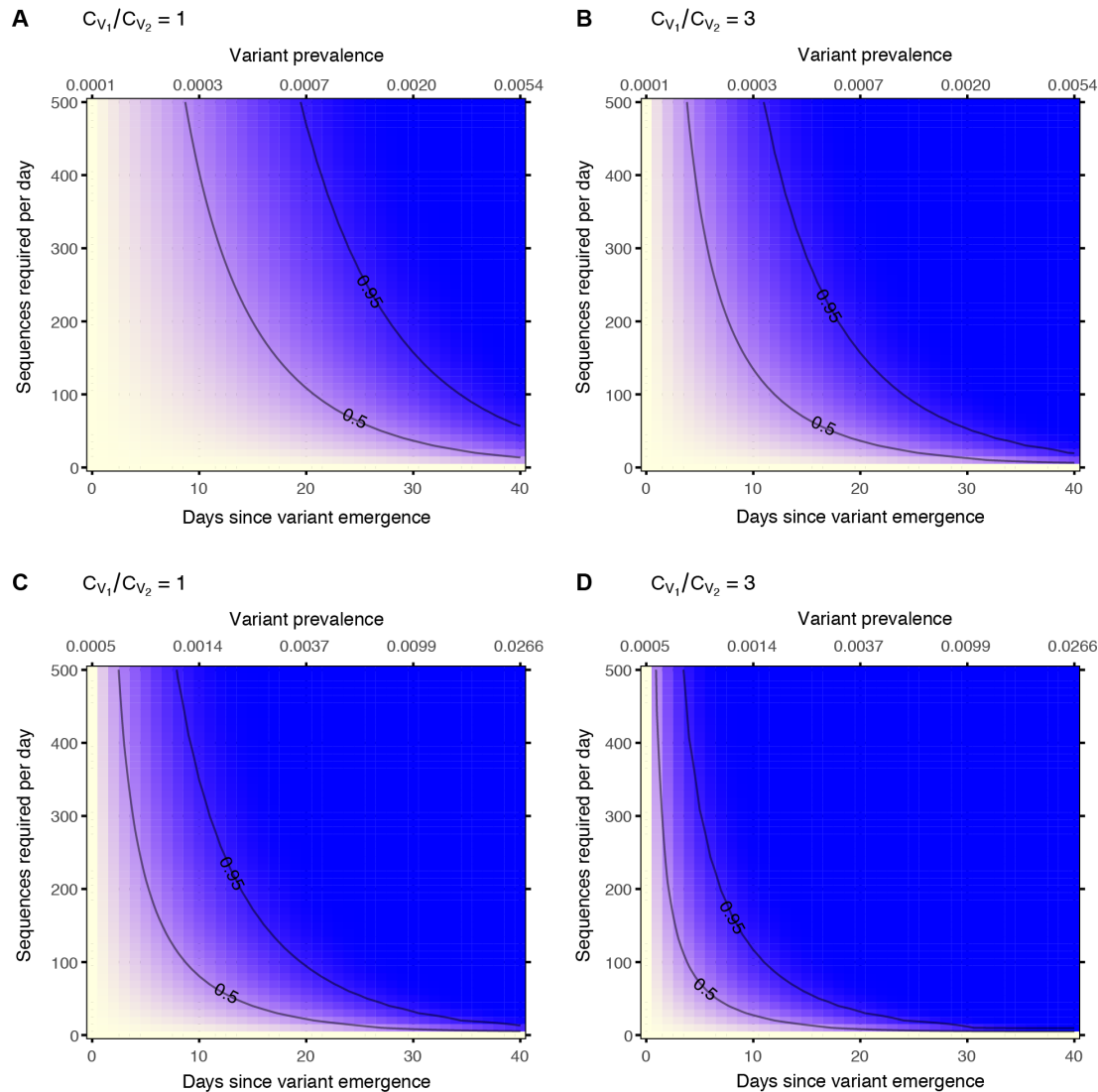


Figure S4. Sample size needed for detection of a variant growing in prevalence. Probability of detecting at least one infection caused by V_1 (yellow = 0% probability; blue = 100% probability) on or before a specific day (bottom x-axis) or desired prevalence (top x-axis), given per-day sample size, specified coefficient of detection ratio, and the following growth rate and initial variant prevalence values: **(A-B)** Growth rate = 0.05, initial prevalence = 1/10000. **(C-D)** Growth rate = 0.1, initial prevalence = 5/10000. Note that the desired prevalence (top x-axis) is the actual variant prevalence in the population and that the number of samples selected for sequencing should exceed the number of sequences required if $\omega < 1$. 50% and 95% probability of detection contours are indicated on all panels.

Appendix

A Calculating variant prevalences from an observed sample

Given the coefficient of detection (C_{V_i}) for each variant in the pathogen population, we can calculate the actual prevalence of each variant (P_{V_i}) from what we observed in the pool of high quality detected infections (H). In the sections below, we use a property of odds ratios to calculate a correction factor q that allows for this conversion.

A.1 Correction factor in a 2-variant system

In a two-variant system (i.e., a system with one variant of interest, V_1 , that is compared to the rest of the population, V_2), the odds of V_1 is:

$$\text{odds}_{V_1} = \frac{P_{V_1}}{P_{V_2}} = \frac{P_{V_1}}{1 - P_{V_1}}$$

Similarly, the observed odds (the odds of P_{V_1} in H) is:

$$\text{odds}_{V_1}^* = \frac{P_{V_1}^*}{P_{V_2}^*} = \frac{\frac{C_{V_1} P_{V_1} N}{C_{V_1} P_{V_1} N + C_{V_2} P_{V_2} N}}{\frac{C_{V_2} P_{V_2} N}{C_{V_1} P_{V_1} N + C_{V_2} P_{V_2} N}} = \frac{C_{V_1} P_{V_1}}{C_{V_2} P_{V_2}} = \frac{C_{V_1} P_{V_1}}{C_{V_2} (1 - P_{V_1})}$$

We define a bias factor, q , such that:

$$\begin{aligned} \text{odds}_{V_1} &= q \times \text{odds}_{V_1}^* \\ \frac{P_{V_1}}{1 - P_{V_1}} &= q \times \frac{C_{V_1} P_{V_1}}{C_{V_2} (1 - P_{V_1})} \end{aligned}$$

If we solve the above equation for q , we obtain:

$$q = \frac{C_{V_2}}{C_{V_1}}$$

Because we know that: $P_{V_i} = \frac{\text{odds}_{V_i}}{1 + \text{odds}_{V_i}}$, we can use the correction factor q to easily calculate the true proportion of any variant i in a population from its observed proportion in the sample H :

$$P_{V_i} = \frac{\text{odds}_{V_i}}{1 + \text{odds}_{V_i}} = \frac{q(\text{odds}_{V_i}^*)}{1 + q(\text{odds}_{V_i}^*)}$$

A.2 Correction factor in a 3-variant system

In a 3-variant system, the odds of V_1 are as follows:

$$\text{odds}_{V_1}^* = \frac{P_{V_1}^*}{P_{V_2}^* + P_{V_3}^*} = \frac{\frac{C_{V_1} P_{V_1} N}{C_{V_1} P_{V_1} N + C_{V_2} P_{V_2} N + C_{V_3} P_{V_3} N}}{\frac{C_{V_2} P_{V_2} N + C_{V_3} P_{V_3} N}{C_{V_1} P_{V_1} N + C_{V_2} P_{V_2} N + C_{V_3} P_{V_3} N}} = \frac{C_{V_1} P_{V_1}}{C_{V_2} P_{V_2} + C_{V_3} P_{V_3}}$$

Using this, we can calculate $q_{V_1,123}$, the correction factor between the true and observed odds when V_1 is the variant of interest in a 3-variant system with V_1 , V_2 , and V_3 :

$$\begin{aligned}
 \frac{P_{V_1}}{1 - P_{V_1}} &= q \times \frac{C_{V_1} P_{V_1}}{C_{V_2} P_{V_2} + C_{V_3} P_{V_3}} = q \times \frac{C_{V_1} P_{V_1}}{(C_{V_2} P_{V_2} + C_{V_3} P_{V_3}) \left(\frac{1 - P_{V_1}}{1 - P_{V_1}} \right)} \\
 &= q \left(\frac{P_{V_1}}{1 - P_{V_1}} \right) \frac{C_{V_1}}{\frac{C_{V_2} P_{V_2}}{1 - P_{V_1}} + \frac{C_{V_3} P_{V_3}}{1 - P_{V_1}}} \\
 &= q \left(\frac{P_{V_1}}{1 - P_{V_1}} \right) \frac{C_{V_1}}{\frac{C_{V_2} P_{V_2}}{P_{V_2} + P_{V_3}} + \frac{C_{V_3} P_{V_3}}{P_{V_2} + P_{V_3}}} \\
 &= q \left(\frac{P_{V_1}}{1 - P_{V_1}} \right) \frac{C_{V_1}}{C_{V_2} \left(\frac{P_{V_2}}{P_{V_2} + P_{V_3}} \right) + C_{V_3} \left(\frac{P_{V_3}}{P_{V_2} + P_{V_3}} \right)}
 \end{aligned}$$

At this point, we recognize that $\frac{P_{V_2}}{P_{V_2} + P_{V_3}}$ is equivalent to P_{V_2} if variants 2 and 3 are the only variants in that system. Assuming a 2-variant system with variants 2 and 3 only, we can use the results of the previous section to write P_{V_2} as a function of the observed odds in this system:

$$P_{V_2} = \frac{q_{V_{2,23}} \text{odds}_{V_{2,23}}^*}{1 + q_{V_{2,23}} \text{odds}_{V_{2,23}}^*}$$

where $\text{odds}_{V_{2,23}}^*$ is the observed odds of V_2 in this 2-variant system with V_2 and V_3 (i.e., $\frac{P_{V_2}^*}{P_{V_3}^*}$) and $q_{V_{2,23}}$ is the correction factor in this system (which we know to be equal to $\frac{C_{V_3}}{C_{V_2}}$). Therefore, we can continue our calculation of $q_{V_{1,123}}$ by substituting these values as follows:

$$\begin{aligned}
 \frac{P_{V_1}}{1 - P_{V_1}} &= q \left(\frac{P_{V_1}}{1 - P_{V_1}} \right) \frac{C_{V_1}}{C_{V_2} \frac{q_{V_{2,23}} \text{odds}_{V_{2,23}}^*}{1 + q_{V_{2,23}} \text{odds}_{V_{2,23}}^*} + C_{V_3} \frac{q_{V_{3,23}} \text{odds}_{V_{3,23}}^*}{1 + q_{V_{3,23}} \text{odds}_{V_{3,23}}^*}} \\
 q &= \frac{1}{C_{V_1}} \left(C_{V_2} \frac{q_{V_{2,23}} \text{odds}_{V_{2,23}}^*}{1 + q_{V_{2,23}} \text{odds}_{V_{2,23}}^*} + C_{V_3} \frac{q_{V_{3,23}} \text{odds}_{V_{3,23}}^*}{1 + q_{V_{3,23}} \text{odds}_{V_{3,23}}^*} \right)
 \end{aligned}$$

A.3 Correction factor in an n-variant system

We can extend the conclusion above to derive a formula for the correction factor q in a system with n variants:

$$q_{V_{1,12..n}} = \frac{1}{C_{V_1}} \left(C_n \frac{q_{V_{n,2..n}} (\text{odds}_{V_{n,2..n}}^*)}{1 + q_{V_{n,2..n}} (\text{odds}_{V_{n,2..n}}^*)} + C_{n-1} \frac{q_{V_{n-1,2..n}} (\text{odds}_{V_{n-1,2..n}}^*)}{1 + q_{V_{n-1,2..n}} (\text{odds}_{V_{n-1,2..n}}^*)} + \dots + C_2 \frac{q_{V_{2,2..n}} (\text{odds}_{V_{2,2..n}}^*)}{1 + q_{V_{2,2..n}} (\text{odds}_{V_{2,2..n}}^*)} \right)$$

The exact value of q , and therefore the true value of P_{V_1} , can be calculated recursively given only the coefficients of detection and observed proportions of the variants in the population.

B Sample size calculations with continuous surveillance

B.1 Sample size calculation for variant detection on or before time t

The probability of detecting a variant (i.e., generating one or more high quality sequences indicating a patient was infected by this variant) on or before time t is equal to one minus the probability of not detecting it at any time between t_0 and t . In other words, regardless of the time unit used, this probability can be written as:

$$\Pr(d \leq t) = 1 - \prod_{x=0}^t \left[1 - \Pr(\text{detection at time } x) \right]$$

Assuming a binomial sampling process, the probability of detection at time x is equal to one minus the probability of not detecting the variant:

$$\Pr(\text{detection at time } x) = 1 - (1 - P_x)^n$$

Where n is the sample size and P_x is the prevalence of the variant at this time step. Therefore, we can write the probability of detecting the variant on or before time t as:

$$\Pr(d \leq t) = 1 - \prod_{x=0}^{x=t} \left[1 - (1 - (1 - P_x)^n) \right] = 1 - \prod_{x=0}^{x=t} (1 - P_x)^n$$

Where n is the per-time step sample size and P_x is the prevalence of the variant at time x . This assumes the same number of samples are selected at every time step, and that the prevalence of the variant at each time step is known. We can rewrite this equation to solve for the per-time step sample size:

$$n = \frac{\ln[1 - \Pr(d \leq t)]}{\ln[\prod_{x=0}^t (1 - P_x)]}$$

We can approximate the value of the product with a continuous function using the Volterra product integral. For a scalar function f and real values of a and b :

$$\prod_a^b (1 + f(x)dx) = \lim_{\Delta x \rightarrow 0} \prod (1 + f(x_i)\Delta x) = \exp\left(\int_a^b f(x)dx\right)$$

Let $f(x) = -P_x dx$. This allows us to write the product of $1 - P_x$ as:

$$\prod_{x=0}^t (1 + (-P_x dx)) = \exp\left(\int_0^t -P_x dx\right) = \frac{1}{\exp(\int_0^t P_x dx)}$$

If we plug this into our per-time step sample size calculation we get:

$$n = \frac{\ln[1 - \Pr(d \leq t)]}{\ln\left[\frac{1}{\exp(\int_0^t P_x dx)}\right]} = \frac{\ln[1 - \Pr(d \leq t)]}{\ln(1) - \ln[\exp(\int_0^t P_x dx)]} = -\frac{\ln[1 - \Pr(d \leq t)]}{\int_0^t P_x dx} = -\frac{\ln[1 - \Pr(d \leq t)]}{G(t) - G(0)}$$

Where $G(t)$ is the integral of $g(t)$. In other words, $G(t)$ is the cumulative density function of the growth function ($g(t) = P_x dx$) used to model the change in the variant frequency over time.

B.2 Logistic growth of variant frequency

Growth in prevalence of variants of interest (i.e., variants with some fitness advantage) are often modeled by logistic growth functions. In other words:

$$g(t) = \frac{1}{1 + ae^{-rt}}$$

Where r is the per-time step growth rate and $a = \frac{1}{t_0} - 1$. The cumulative density of this probability distribution can be computed as follows:

$$\begin{aligned} G(t) &= \int \frac{1}{1 + ae^{-rt}} dt = \int \left(\frac{e^{rt}}{e^{rt}}\right) \frac{1}{1 + ae^{-rt}} dt = \int \frac{e^{rt}}{e^{rt} + a} dt \\ &= \int \frac{1}{u} \frac{du}{r} = \frac{1}{r} \int \frac{1}{u} du = \frac{1}{r} \ln |u| \\ &= \frac{1}{r} \ln |a + e^{rt}| + C \end{aligned}$$

Where $u = e^{rt} + a$.

B.3 Logistic growth of variant frequency with biased detection

In a two-variant system, the observed prevalence of a particular variant of interest in a sample of high quality detected infections (H) is a function of the true prevalence and the ratio between coefficients of detection:

$$\text{observed frequency} = \frac{P_{V_1}}{P_{V_1} + \frac{C_{V_2}}{C_{V_1}} P_{V_2}}$$

Assuming P_{V_1} can be computed at any given time step using a logistic model, we can calculate the observed frequency distribution as follows:

$$\text{observed frequency} = \frac{\frac{1}{1+ae^{-rt}}}{\frac{1}{1+ae^{-rt}} + \frac{C_{V_2}}{C_{V_1}} \left(1 - \frac{1}{1+ae^{-rt}}\right)} = \frac{1}{1 + \frac{C_{V_2}}{C_{V_1}} ae^{-rt}} = \frac{1}{1 + be^{-rt}}$$

Where $b = a \frac{C_{V_2}}{C_{V_1}} = \left(\frac{C_{V_2}}{C_{V_1}}\right) \left(\frac{1}{t_0} - 1\right)$.

Because in this case the observed frequency function takes the same form as the actual frequency function, we can easily calculate $G^*(t)$, the cumulative density of the observed variant frequency:

$$G^*(t) = \frac{1}{r} \ln |b + e^{rt}| + C$$