

## Supplement A

### Variables

- Exposure to SARS-CoV-2 (*Exposure*)
  - Binary. “1” for exposed and “0” for unexposed.
- SARS-CoV-2 Infection (*Infection*)
  - Binary. “1” for infected and “0” for uninfected.
- COVID-19 associated hospitalization (*Hospitalization*)
  - Binary. “1” for hospitalized and “0” for non-hospitalized/mild symptoms.
- Received testing for SARS-CoV-2 infection or not (*If\_tested*)
  - Binary. “1” if the person received testing and “0” otherwise.
- Test result for SARS-CoV-2 infection (*Test\_result*)
  - Binary. “1” if tested positive and “0” if tested negative.
- Number of protection alleles carried: Variant associated with infection susceptibility
  - $G_{inf}$ ; 0, 1, or 2
- Number of risk alleles carried: Variant associated with disease severity / hospitalization
  - $G_{hosp}$ ; 0, 1, or 2

### Parameters

- Population Exposure Rate ( $p_{exposure}$ )
- Minor allele frequency of  $G_{inf}$  ( $MAF_{inf}$ )
- Minor allele frequency of  $G_{hosp}$  ( $MAF_{hosp}$ )
- Baseline Infection Susceptibility ( $p_{suscep}$ )
  - Probability of infection on exposure to SARS-CoV-2 in the absence of the contributing protective genetic allele.
- Baseline risk of hospitalization ( $p_{hosp}$ )
  - Probability of developing severe symptoms/hospitalization given SARS-CoV-2 infection in the absence of the contributing risk allele.
- Test prevalence for individuals with SARS-CoV-2 infection and mild symptoms
  - $p_{test\_prev\_mild}$
  - Test prevalence for individuals with severe symptoms is assumed to be 100%
- Test prevalence for individuals without SARS-CoV-2 infection
  - $p_{test\_prev\_no\_inf}$
- Decreased risk of infection with an additional protection allele of  $G_{inf}$ 
  - $OR_{inf}$ ; Log-additive.
- Increased risk of hospitalization with an additional risk allele of  $G_{hosp}$ 
  - $OR_{hosp}$ ; Log-additive
- Test sensitivity (*sens*)
- Test specificity (*spec*)

## Simulated Data Generating Mechanism

1. **Simulate genetic variants.** Each individual  $i$  carries
  - a. Variant (Infection Susceptibility):  $G_{inf,i} \sim Binom(2, MAF_{inf})$ .
  - b. Variant (Disease Severity / Hospitalization):  $G_{hosp,i} \sim Binom(2, MAF_{hosp})$
2. **Simulate exposure to SARS-CoV-2.**
  - a.  $Exposure_i \sim Bernoulli(p_{exposure})$
3. **Simulate SARS-CoV-2 infection.**
  - a. If  $Exposure_i = 0$ ,  $Infection_i = 0$ .
  - b. If  $Exposure_i = 1$  and  $G_{inf,i} = 0$ ,  $Infection_i \sim Bernoulli(p_{suscep})$
  - c. If  $Exposure_i = 1$  and  $G_{inf,i} \neq 0$ ,  $Infection_i \sim Bernoulli\left(f\left(f^{-1}(p_{suscep}) + G_{inf,i} \times OR_{inf}\right)\right)$ , where  $f$  is the logistic function.
4. **Simulate severe symptoms / hospitalization given SARS-CoV-2 infection.**
  - a. If  $Infection_i = 0$ ,  $Hospitalization_i = 0$ .
  - b. If  $Infection_i = 1$  and  $G_{hosp,i} = 0$ ,  $Hospitalization_i \sim Bernoulli(p_{hosp})$
  - c. If  $Infection_i = 1$  and  $G_{hosp,i} \neq 0$ ,  $Hospitalization_i \sim Bernoulli\left(f\left(f^{-1}(p_{hosp}) + G_{hosp,i} \times OR_{hosp}\right)\right)$  where  $f$  is the logistic function.
5. **Simulate whether the individual received testing for SARS-CoV-2 infection.**
  - a. If  $Infection_i = 0$ ,  $If\_tested_i \sim Bernoulli(p_{test\_prev\_no\_inf})$
  - b. If  $Infection_i = 1$  and  $Hospitalization_i = 0$ ,  $If\_tested_i \sim Bernoulli(p_{test\_prev\_mild})$
  - c. If  $Hospitalization_i = 1$ ,  $If\_tested_i = 1$
6. **Simulate Test Result for SARS-CoV-2 infection.**
  - a. If  $Infection_i = 0$ ,  $Test\_result_i \sim Bernoulli(1 - spec)$
  - b. If  $Infection_i = 1$  and  $Hospitalization_i = 0$ ,  $Test\_result_i \sim Bernoulli(sens)$
  - c. If  $Hospitalization_i = 1$ ,  $Test\_result_i = 1$

## Infection Susceptibility Study Design

### **Inclusion Criteria**

1. Those who received testing for SARS-CoV-2 infection ( $I_{f\_tested_i} = 1$ ).

### **Case-control definition**

1. Cases: Test-positive individuals ( $Test\_result_i = 1$ ).
2. Controls: Test-negative individuals ( $Test\_result_i = 0$ ).

Standard univariate logistic regression was used to estimate the effect size of the genetic variant. A finding is reported if the p-value is below the genome-wide significance threshold,  $5e-8$ .

### Default Parameter Settings

1. Baseline Infection Susceptibility ( $p_{suscep}$ )
  - 80% given current estimates for SARS-CoV-2 infection [11].
2. Baseline risk of hospitalization ( $p_{hosp}$ )
  - 5%. States/Provinces across U.S. and Canada currently report hospitalization rates between 8~10% but may overestimate the figure due to many asymptomatic patients / false negatives.
3. Test prevalence for individuals with SARS-CoV-2 infection and mild symptoms
  - $p_{test\_prev\_mild}$
  - 30%
4. Test prevalence for individuals without SARS-CoV-2 infection
  - $p_{test\_prev\_no\_inf}$
  - 5%

### Case-Control Misclassification

Proportion of population received testing and tested positive for SARS-CoV-2 infection

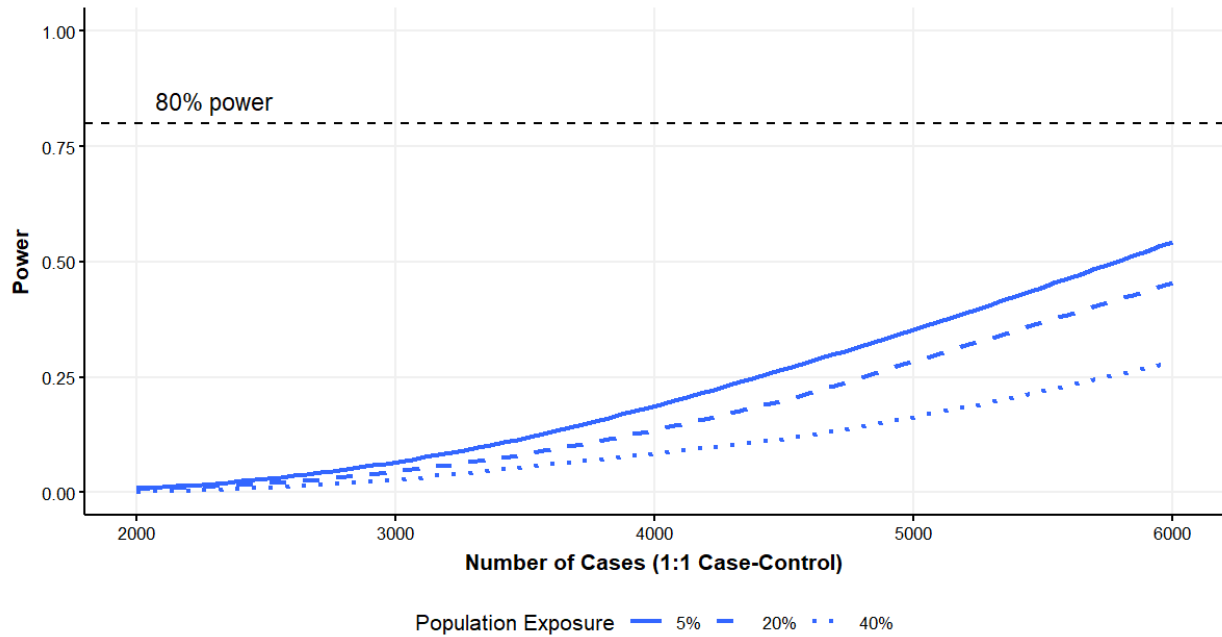
$$total\_cases = p_{exposure} \times p_{suscep} \times (p_{hosp} \times 1 + (1 - p_{hosp}) \times p_{test\_prev\_mild}) \times sens \\ + ((1 - p_{exposure}) + p_{exposure} \times (1 - p_{suscep})) \times p_{test\_prev\_no\_inf} \times (1 - spec)$$

Proportion of population received testing and tested negative for SARS-CoV-2 infection

$$total\_controls = p_{exposure} \times p_{suscep} \times (p_{hosp} \times 1 + (1 - p_{hosp}) \times p_{test\_prev\_mild}) \times (1 - sens) \\ + ((1 - p_{exposure}) + p_{exposure} \times (1 - p_{suscep})) \times p_{test\_prev\_no\_inf} \times spec$$

<b>Misclassification among cases</b>	
True controls (uninfected) misclassified as cases due to false positives produced by RT-PCR tests	$\frac{(1 - p_{exposure}) \times p_{test\_prev\_no\_inf} \times (1 - spec)}{total\_cases}$
<b>Misclassification among controls</b>	
Unexposed individuals that should have been infected upon exposure	$\frac{(1 - p_{exposure}) \times p_{test\_prev\_no\_inf} \times p_{suscep}}{total\_controls}$
True cases (infected) misclassified as controls (uninfected) due to a low sensitivity RT-PCR test	$\frac{p_{exposure} \times p_{suscep} \times (1 - p_{hosp}) \times p_{test\_prev\_mild} \times (1 - sens)}{total\_controls}$

## Supplement B



**Supplementary Figure B: Statistical power to detect associations between genetic variants and infection susceptibility at the genome-wide significance level ( $5e-8$ ) when the test sensitivity is low (sensitivity = 0.7). A 1:1 case-control study design was used for all parameter settings. Assumes a common variant with large effect size (OR=0.5, MAF=0.2). Reducing sensitivity for testing SARS-CoV-2 infection not only reduces statistical power but also negates gains that result from increasing population exposure. In fact, increasing population exposure leads to reduced statistical power when the test sensitivity is too low.**

## Supplement C

Variables, parameters and the simulation process are all identical to those provided in Supplement A.

### Disease Severity Study Design

#### Using Test-Positive Controls

##### Inclusion Criteria

1. Those who received testing for SARS-CoV-2 infection and tested positive ( $If\_tested_i = 1$  &  $test\_result_i = 1$ ).

##### Case-control definition

1. Cases: Individuals with severe symptoms (hospitalized;  $Hospitalization_i = 1$ ).
2. Controls: Individuals without severe symptoms ( $Hospitalization_i = 0$ ).

#### Using Population-based (untested) Controls

##### Inclusion Criteria

1. Those who received SARS-CoV-2 infection, tested positive and were hospitalized ( $Hospitalization_i = 1$ ), and the population-based controls that never received testing ( $If\_tested_i = 0$ ).

##### Case-control definition

1. Cases: Individuals with severe symptoms (hospitalized;  $Hospitalization_i = 1$ ).
2. Controls: Population-based controls that never received testing ( $If\_tested_i = 0$ ).

Standard univariate logistic regression was used to estimate the effect size of the genetic variant. A finding is reported if the p-value is below the genome-wide significance threshold,  $5e-8$ .

### Default Parameter Settings

Contrary to the low population exposure that drives case-control misclassification in SARS-CoV-2 infection susceptibility studies, it is the low population *infection* rates and the lack of information on individual-level *infection* that can lead to misclassification in disease severity studies; as some of the defined “controls” would have developed severe symptoms had they been infected with SARS-CoV-2.

1. Baseline Infection Susceptibility ( $p_{suscep}$ )

- 100%. Since it is the low population *infection* rates and the lack of information on individual-level *infection* that can lead to misclassification in disease severity studies. Varying population exposure is equivalent to varying population infection rates if  $p_{suscep} = 1$ .
2. Baseline risk of hospitalization ( $p_{hosp}$ )
    - 5%. States/Provinces across U.S. and Canada currently report hospitalization rates between 8~10% but may overestimate the figure due to many asymptomatic patients / false negatives.
  3. Test prevalence for individuals with SARS-CoV-2 infection and mild symptoms
    - $p_{test\_prev\_mild}$
    - 30%
  4. Test prevalence for individuals without SARS-CoV-2 infection
    - $p_{test\_prev\_no\_inf}$
    - 5%.

### Case-Control Misclassification

Define population infection rate as  $p_{inf} = p_{exposure} \times p_{suscep}$ .

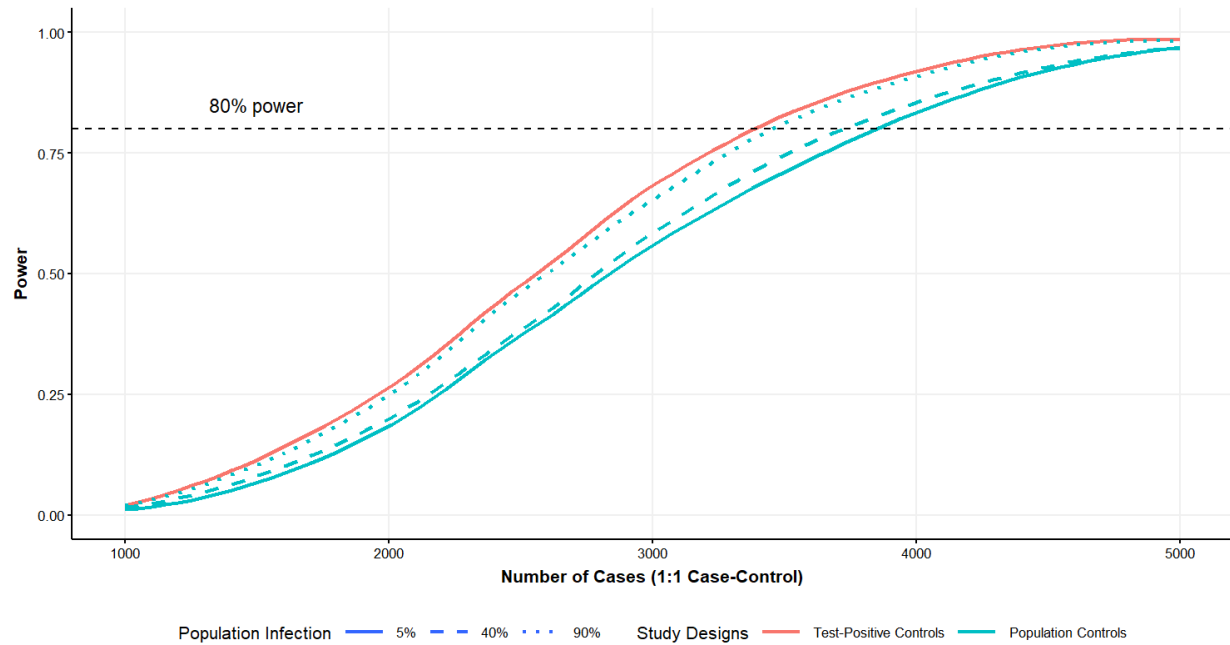
#### **Using population-based (uninfected) controls**

Misclassification of controls	
Uninfected individuals that should have been hospitalized upon infection	$(1 - p_{inf}) \times p_{hosp}$

#### **Using test-positive controls**

Misclassification of controls	
Test-positive controls (infected) which were in fact not infected due to false positives produced by RT-PCR tests. Some would have developed severe symptoms upon infection.	$\frac{(1 - p_{inf}) \times (1 - spec)}{p_{inf} \times sens + (1 - p_{inf}) \times (1 - spec)} \times p_{hosp}$

## Supplement D



**Supplementary Figure D: Statistical power to detect a true association between a genetic variant and COVID-19 disease severity at the genome-wide significance level ( $5e-8$ ). A 1:1 case-control study design was used for all parameter settings. Only one red curve is shown since the study design uses confirmed infected individuals with mild or no symptoms as controls (test-positive controls), which is unaffected by population-level infection rates and the corresponding case-control misclassification. Effect sizes are reported on the odds ratio (OR) scale for each additional risk allele. Assumes perfect test accuracy. Detecting a common variant with moderate effect size (OR=1.3, MAF=0.2) is much more challenging without drastically increasing the number of participants studied. In fact, the sample size required to reach 80% power exceeds 4 times that required when the effect size is large (OR=1.7).**

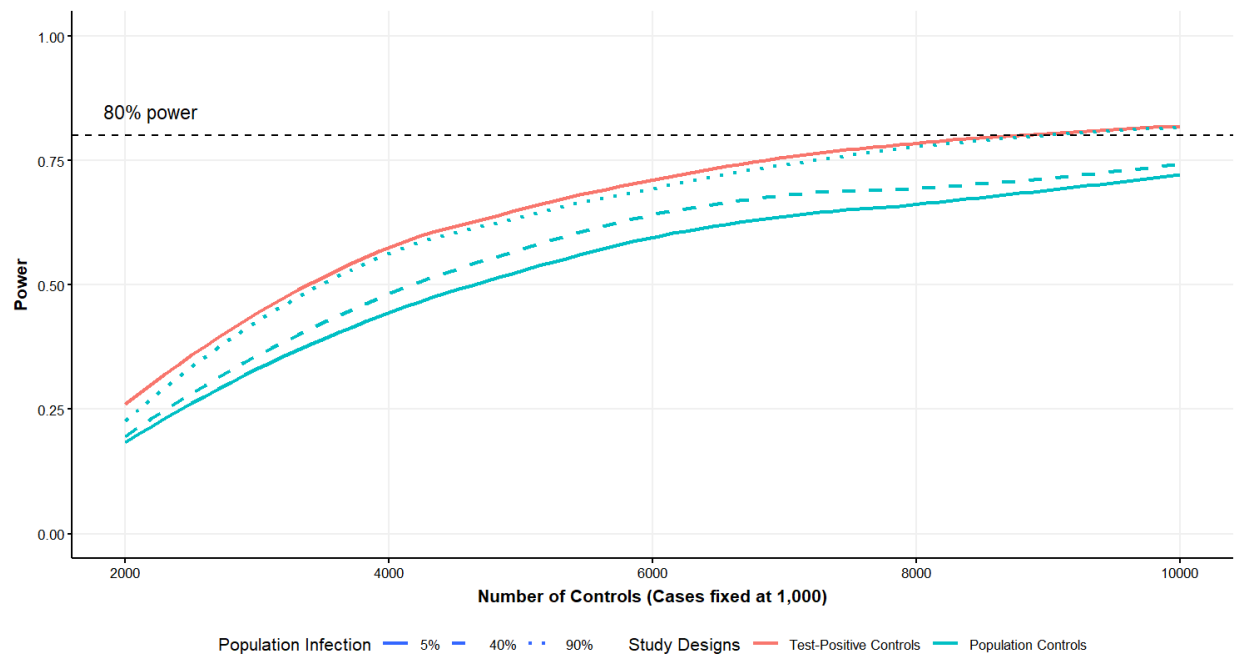


## Supplement E

Sample Size Required to Achieve 80% Power in Detecting True Associations (1:1 Case-Control Setting)			
Parameter Settings / Study Designs	Test-Positive Controls	Population-based (Untested) Controls	Gain in Efficiency ( $\frac{\text{Sample Size}_{B2}}{\text{Sample Size}_{B1}}$ )
<b>Large Effect Size (OR=1.7; MAF=0.2)</b>			
<i>Population Infection Rates: 0.1</i>	<b>800</b>	910	13.8%
<i>0.2</i>		900	12.5%
<i>0.3</i>		890	11.2%
<i>0.4</i>		880	10.0%
<i>0.5</i>		870	8.7%
<i>0.6</i>		855	6.9%
<i>0.7</i>		845	5.6%
<i>0.8</i>		830	3.7%
<i>0.9</i>		815	1.9%
<b>Moderate Effect Size (OR=1.3; MAF=0.2)</b>			
<i>Population Infection Rates: 0.1</i>	<b>3,390</b>	3,800	12.1%
<i>0.2</i>		3,760	10.9%
<i>0.3</i>		3,720	9.7%
<i>0.4</i>		3,680	8.6%
<i>0.5</i>		3,640	7.4%
<i>0.6</i>		3,595	6.0%
<i>0.7</i>		3,550	4.7%
<i>0.8</i>		3,500	3.2%
<i>0.9</i>		3,445	1.6%
<b>Rare Variant with Very Large Effect Sizes (OR=5; MAF=0.01)</b>			
<i>Population Infection Rates: 0.1</i>	<b>1,025</b>	1,150	12.2%
<i>0.2</i>		1,135	10.7%
<i>0.3</i>		1,120	9.3%
<i>0.4</i>		1,110	8.3%
<i>0.5</i>		1,100	7.3%
<i>0.6</i>		1,090	6.3%
<i>0.7</i>		1,075	4.9%
<i>0.8</i>		1,060	3.4%
<i>0.9</i>		1,045	1.9%

**Supplementary Table E: Relative reduction in sample size,  $1 - \frac{n_{test\_positive\_controls}}{n_{population\_controls}}$ , from using test-positive controls compared to population-based controls.  $n_{test\_positive\_controls}$  and  $n_{population\_controls}$  refer to the number of cases (1:1 case-control ratio) needed to achieve 80% power at the genome-wide significance level ( $5e-8$ ) [22]. Using controls with confirmed infection show the greatest benefit when disease prevalence is low. Given current levels of population infection rates, choosing controls with confirmed SARS-CoV-2 infection can save over 10% in genotyping costs.**

## Supplement F



**Supplementary Figure F: Statistical power to detect a true association between a genetic variant and COVID-19 disease severity at the genome-wide significance level ( $5e-8$ ). Cases are fixed at 1,000 and controls are varied to allow reduced case-control ratio. Effect sizes are reported on the odds ratio (OR) scale for each additional risk allele. Assumes perfect test accuracy. Detecting a common variant with moderate effect size (OR=1.3, MAF=0.2) requires much larger sample sizes if a reduced case-control ratio is used.**