



INSTITUTE FOR DEFENSE ANALYSES

**Understanding Viral Population Genomics:
Tools and Lessons for Future Pandemics**
(Poster)

Katherine I. Fisher-Aylor
Izzy Chaiken
Emily D. Heuring

September 2022

Approved for public release;
distribution is unlimited.

IDA Document NS D-33249

Log: H 22-000397

INSTITUTE FOR DEFENSE ANALYSES
730 East Glebe Road
Alexandria, Virginia 22305-3086



The Institute for Defense Analyses is a nonprofit corporation that operates three Federally Funded Research and Development Centers. Its mission is to answer the most challenging U.S. security and science policy questions with objective analysis, leveraging extraordinary scientific, technical, and analytic expertise.

About This Publication

This work was conducted by the Institute for Defense Analyses Central Research Program, project C2277 “SARS-CoV-2 Paper.” The views, opinions, and findings should not be construed as representing the official position of either the Department of Defense or the sponsoring organization.

For More Information

Katherine I. Fisher-Aylor, Project Leader
kfisher@ida.org, 703-845-6902

Leonard J. Buckley, Director, Science and Technology Division
lbuckley@ida.org, 703-578-2800

Copyright Notice

© 2022 Institute for Defense Analyses
730 East Glebe Road, Alexandria, Virginia 22305-3086 • (703) 845-2000.

This material may be reproduced by or for the U.S. Government pursuant to the copyright license under the clause at DFARS 252.227-7013 (Feb. 2014).

The Covid-19 pandemic is the most information-rich pandemic in history.

For the first time, we were able to obtain entire genome sequences of the pathogen in real time during the pandemic.

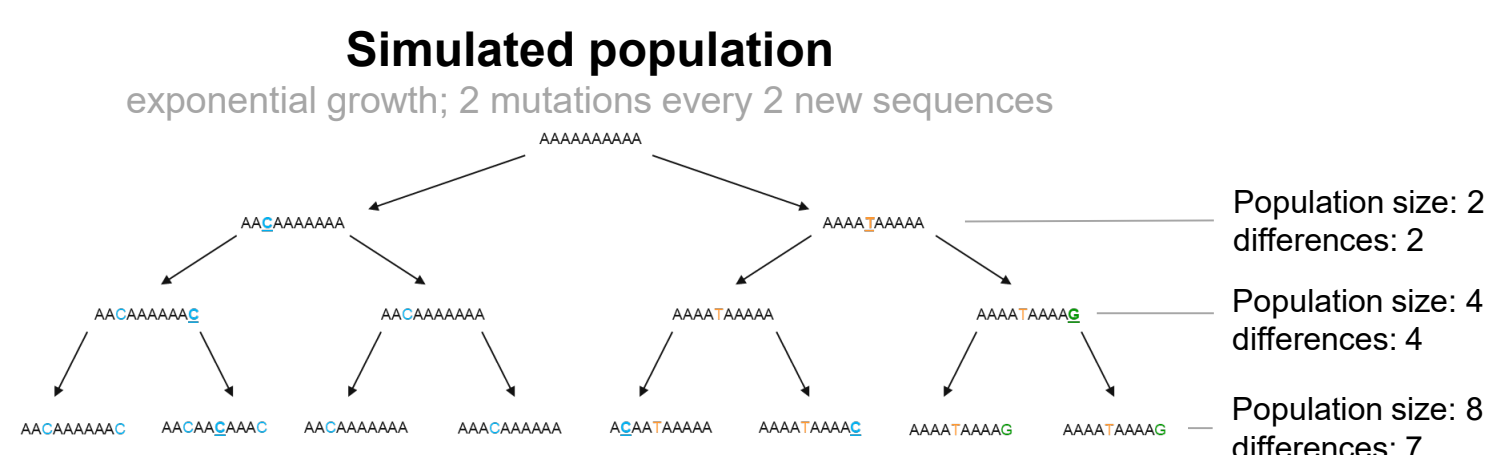
Although sampling of genomes was spotty in places, we had an unprecedented opportunity to understand principles of viral population genomics on an extremely large scale.

During the first year of the pandemic, we wanted to know:

- How many cases are there? Is it possible to use viral genomes to estimate this?
- How much undersampling is too much? Many places in the U.S. collected very few SARS-CoV-2 samples. How much can we undersample before we get a skewed picture of the pandemic?

POPULATION GENETICS PRINCIPLES

Mutation numbers increase with population size.



Therefore, it is possible to infer population size from mutation number.

In the field of population genetics, there are several canonical ways to measure mutations in a population (genetic variation).

Theta
number of sites with a difference
~normalized by population size

$$\theta = \frac{S}{\sum_{i=1}^{n-1} \frac{1}{i}}$$

$$\theta = 2N\mu$$

S: number of segregating sites (places that have a mutation somewhere in the population)
n: number of samples
 μ : mutation rate (for the bounds of the whole region or genome in question)
N: population size
i (left): index of summation (which sample you are on)
j (right): two sequences
 x_i, x_j : frequency of two (i^{th} and j^{th}) sequences
 π_{ij} : number of differences between the two sequences (i and j)

Pi
average pairwise differences

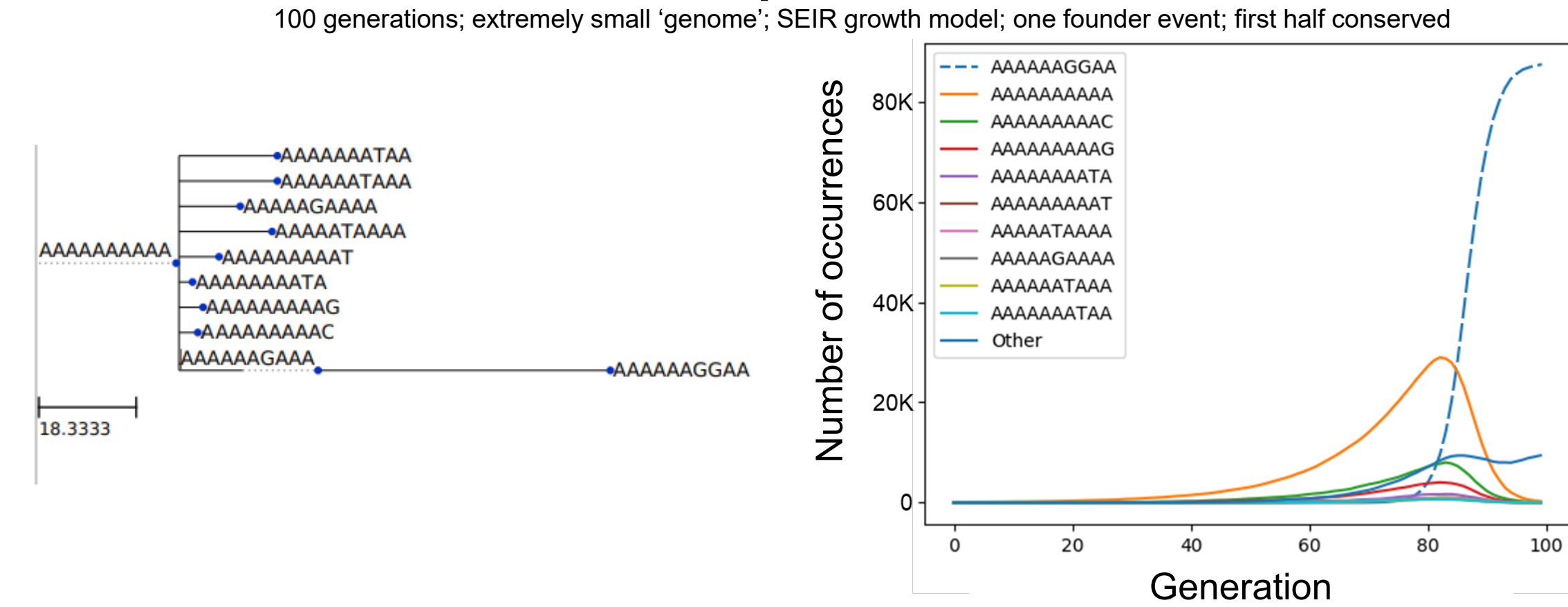
$$\pi = \sum_{i,j} x_i x_j \pi_{ij}$$

$$\pi = \frac{\text{sum of pairwise differences}}{\text{number of pairs}}$$

We analyzed the appropriateness of these canonical estimators for measuring SARS-CoV-2 genetic variation during the 2020 Covid-19 pandemic.

- We created this software package to model a growing and mutating population of DNA sequences.
- Simdemic allows specification of mutation rate, growth model, genome size, selective pressure, and founder events.
- We have used Simdemic to create an explicitly known population from which we can infer how real-world methodology affects our conclusions.
 - Assess different estimators of genetic variation
 - Learn the limits of subsampling
- In the future, we hope to use Simdemic to refine predictions of driving mutations, explore the accuracy of different epidemiological growth models, and understand the impact of founder events (as opposed to selective pressure) on a population.

Example simulation



RESULTS

Theta is an inappropriate estimator of SARS-CoV-2 genetic variation; Pi is appropriate

Simulated data

Simdemic shows theta is the same regardless of number of mutations; pi is not

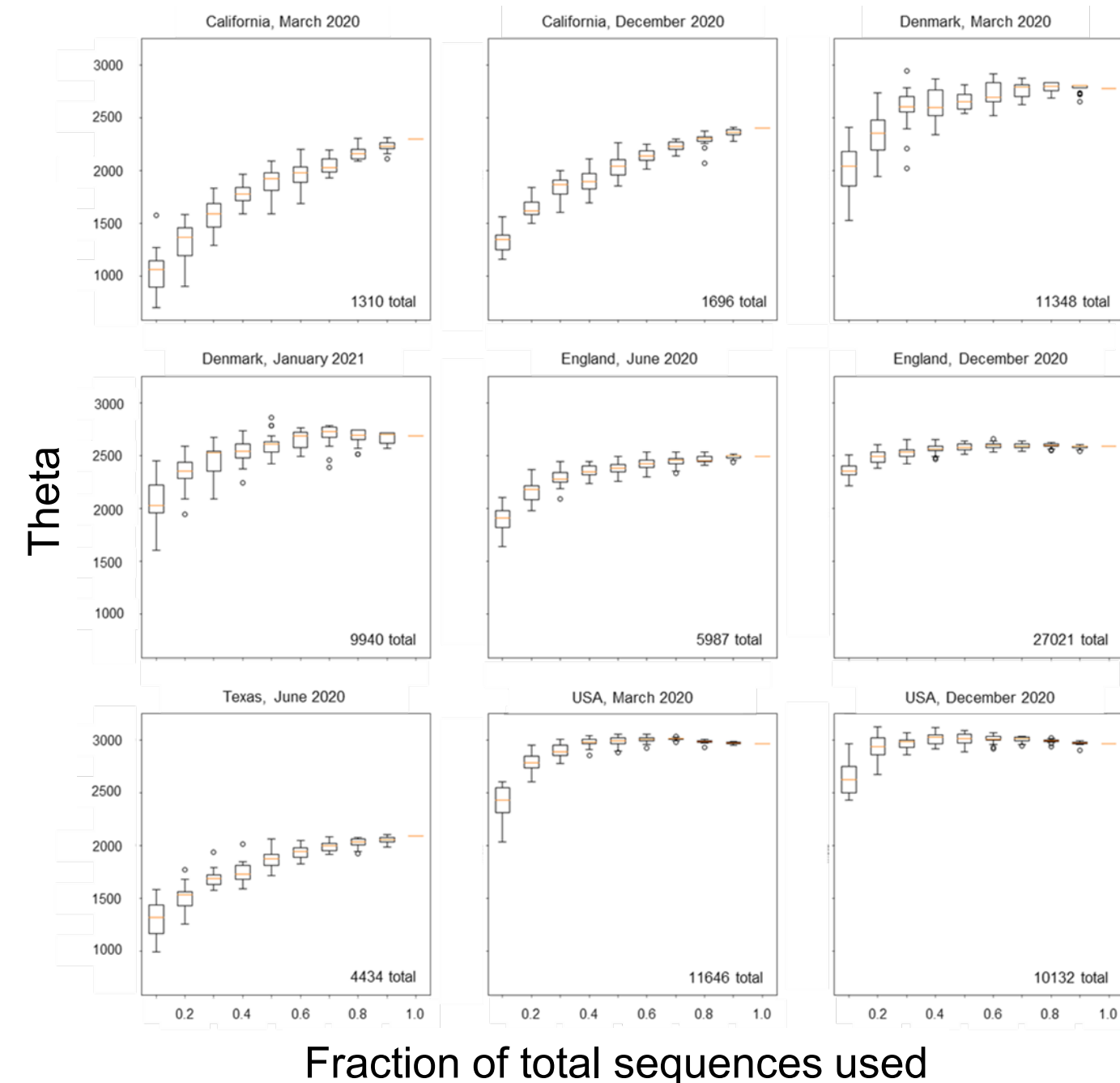
```

sequences 1824
unique 728
[("AACAAAGCAA", 18), ("AAAAAGCAA", 9), ("AAATACTACA", 7), ("AAATACAAA", 7), ("AAAGGCAAA", 6), ("AAGGAATAT", 6),
("AAGAAACAC", 4), ("GAATATACCA", 4), ("AAGAACTACA", 4), ("TAAAGACAAA", 4), ("TATATACCA", 4), ("TACAGCAAA", 4), ("TATACACAAA", 4),
("CAATATACAA", 4), ("AAAAACTAAA", 4), ("AAAAACAG", 4), ("GAAAAACAG", 4), ("AAGAAATA", 4), ("AACACAAA", 4), ("AAAAATTGAG", 4),
("AAGAAAGAG", 4), ("AAGAAAGAG", 4), ("AAAAATACG", 4), ("AAGATCAAA", 3)]
seq length 18
seg sites 18
theta 1.8
pi 0.11177777777777778

```

Real-world data

Theta is insensitive for SARS-CoV-2 populations >1,000



Underlying math

Number of nucleotides: 10
Sites with a difference: 10
Average pairwise difference: 2

This quickly saturates: every nucleotide in a small genome will soon have a difference at every position at least once

$$\theta = \frac{S}{\sum_{i=1}^{n-1} \frac{1}{i}}$$

This takes a long time to saturate

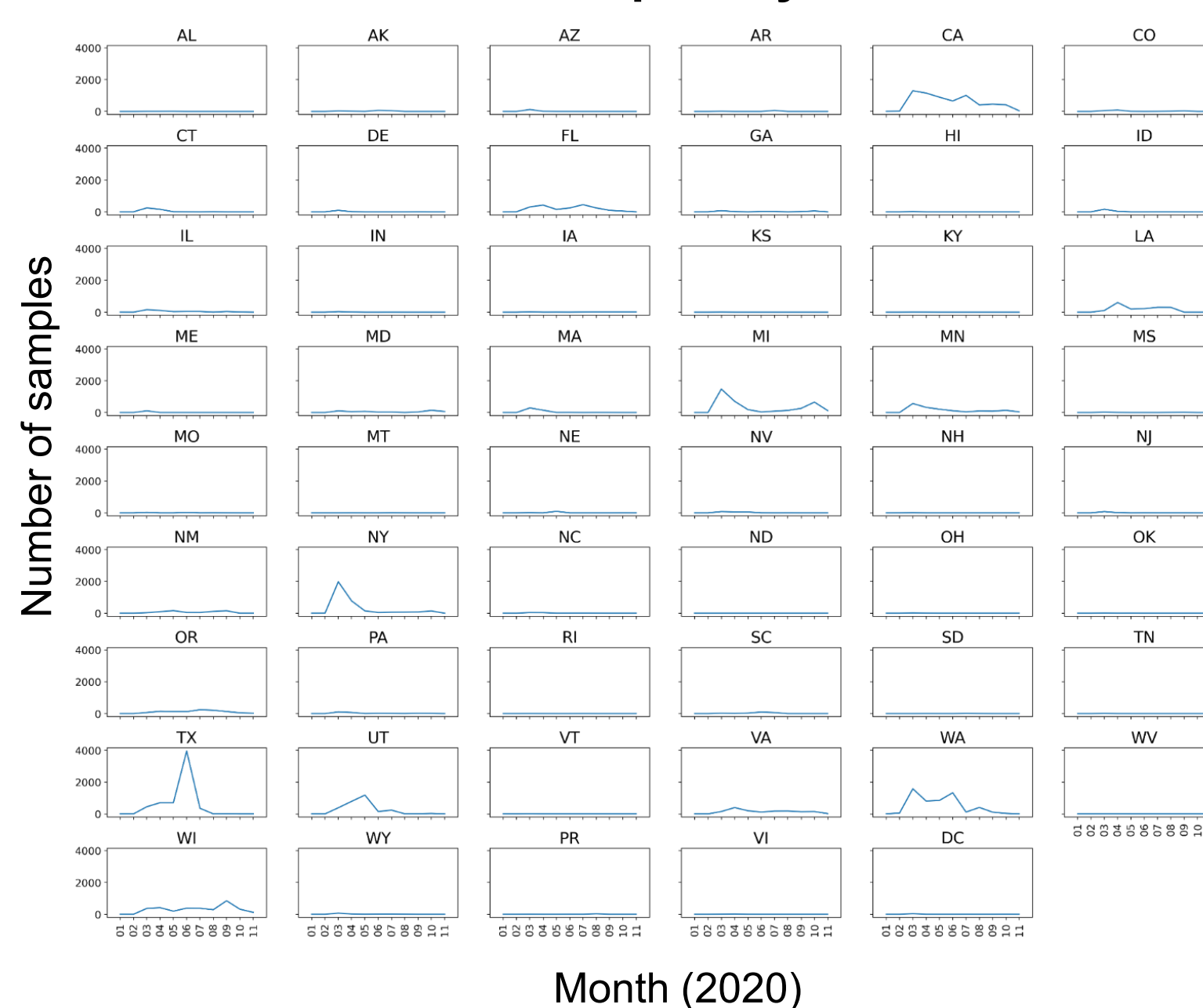
$$\pi = \frac{\text{sum of pairwise differences}}{\text{number of pairs}}$$

Subsampling a population below 1% skews results

GISAID SARS-CoV-2 samples per capita by country Jan-Nov 2020

Country	SARS-CoV-2 sequences	Population size	Ratio
USA	80 K	328 M	1 : 4,100
UK	184 K	67 M	1 : 364
Denmark	32 K	6 M	1 : 188

GISAID SARS-CoV-2 samples by state Jan-Nov 2020



Subsampling with Simdemic

$$\text{Loss function: } \sum_{s \in S} \left(\frac{n_s}{n} - p_s \right)^2$$

S: set of all sequences in the population
 n_s : count of sequence s in the sample
n: sample size
 p_s : true proportion of sequence s in the population

- No centralized effort to collect samples in the U.S.
- Many states have very few samples, making pandemic predictions extremely difficult.
- During the pandemic, we often had to extrapolate from well-sampled countries.

- Preliminary work: 1% is the threshold for subsampling
- Denmark and the UK sampled ~1% of population, while U.S. sampled ~0.01%
- In the future, this can be used to estimate how undersampled a population is, and how many more samples would be necessary to get a fair assessment of a pandemic.

Like other viruses,



Genome

- S
- SA
- E. co
- Br
- F
- C
- Z
- R

C

- Caution must be taken in interpreting canonical population genetic results to small genomes
- Theta is a poor estimator of genetic variation in a population; Pi is an appropriate estimator
- Subsampling a population below 1% skews results, which is a problem for the U.S. rate of sampling during the 2019 pandemic
- More research is needed to understand the relationship between population size and human sampling, like Simdemic, to better understand population genetics in small genomes

REPORT DOCUMENTATION PAGE

*Form Approved
OMB No. 0704-0188*

The public reporting burden for this collection of information is estimated to average 1 hour per response, including the time for reviewing instructions, searching existing data sources, gathering and maintaining the data needed, and completing and reviewing the collection of information. Send comments regarding this burden estimate or any other aspect of this collection of information, including suggestions for reducing the burden, to Department of Defense, Washington Headquarters Services, Directorate for Information Operations and Reports (0704-0188), 1215 Jefferson Davis Highway, Suite 1204, Arlington, VA 22202-4302. Respondents should be aware that notwithstanding any other provision of law, no person shall be subject to any penalty for failing to comply with a collection of information if it does not display a currently valid OMB control number.

PLEASE DO NOT RETURN YOUR FORM TO THE ABOVE ADDRESS.

1. REPORT DATE (DD-MM-YYYY)		2. REPORT TYPE		3. DATES COVERED (From - To)	
4. TITLE AND SUBTITLE				5a. CONTRACT NUMBER	
				5b. GRANT NUMBER	
				5c. PROGRAM ELEMENT NUMBER	
6. AUTHOR(S)				5d. PROJECT NUMBER	
				5e. TASK NUMBER	
				5f. WORK UNIT NUMBER	
7. PERFORMING ORGANIZATION NAME(S) AND ADDRESS(ES)				8. PERFORMING ORGANIZATION REPORT NUMBER	
9. SPONSORING/MONITORING AGENCY NAME(S) AND ADDRESS(ES)				10. SPONSOR/MONITOR'S ACRONYM(S)	
				11. SPONSOR/MONITOR'S REPORT NUMBER(S)	
12. DISTRIBUTION/AVAILABILITY STATEMENT					
13. SUPPLEMENTARY NOTES					
14. ABSTRACT					
15. SUBJECT TERMS					
16. SECURITY CLASSIFICATION OF:			17. LIMITATION OF ABSTRACT	18. NUMBER OF PAGES	19a. NAME OF RESPONSIBLE PERSON
a. REPORT	b. ABSTRACT	c. THIS PAGE			19b. TELEPHONE NUMBER (Include area code)