

# SimPrily:

## A Python framework to simplify genome simulation with priors

Ariella L. Gladstein<sup>1</sup>, Consuelo D. Quinto-Cortés<sup>2</sup>, Julian L. Pistorius<sup>3</sup>, David Christy<sup>4</sup>, Logan Gantner<sup>5</sup>, August E. Woerner<sup>6</sup>, and Blake L. Joyce<sup>3,7</sup>

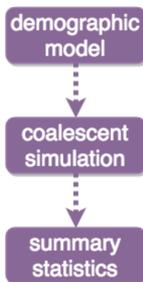
<sup>1</sup>Department of Ecology and Evolutionary Biology, University of Arizona, USA, <sup>2</sup>National Laboratory of Genomics for Biodiversity (LANGEBIO), CINVESTAV, Mexico, <sup>3</sup>CyVerse, University of Arizona, USA, <sup>4</sup>Department of Computer Science, University of Arizona, USA, <sup>5</sup>Graduate Interdisciplinary Program in Applied Mathematics, University of Arizona, USA, <sup>6</sup>Center for Human Identification, University of North Texas Health Science Center, USA, <sup>7</sup>BIO5 Institute, University of Arizona, USA.

**Availability:** Source code, HTC workflow, documentation, and examples are available at <https://github.com/agladstein/SimPrily>

### Introduction

#### What can you use 1000's – millions of simulations for?

- Approximate Bayesian Computation to infer demographic history
- Null demographic model to find regions under selection
- Truth datasets for testing software



#### Features

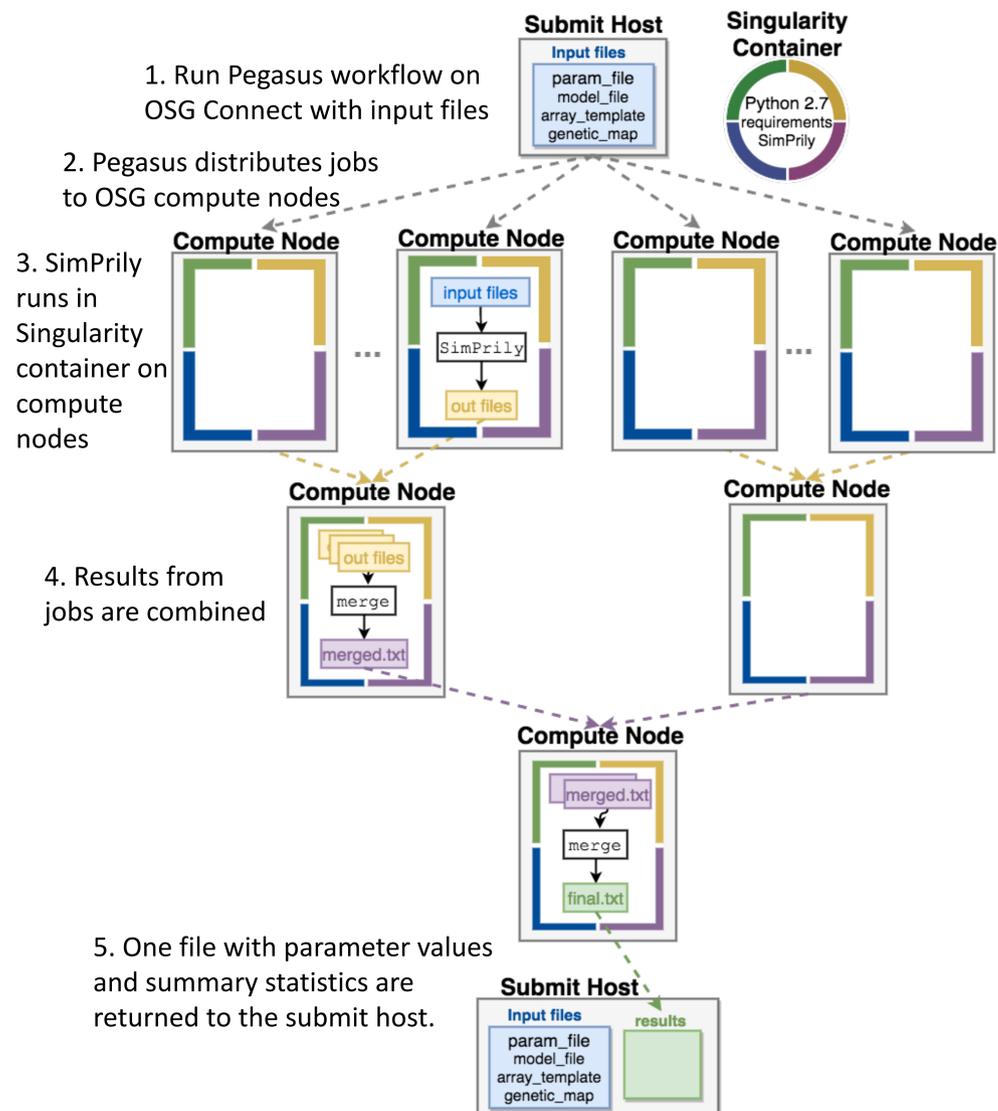
- Specify demographic model with priors
- Create pseudo array from simulations
- Calculate population genetics statistics
- Run 1000's of simulations with GUI in CyVerse Discovery Environment
- Run millions of simulations with Pegasus workflow on the Open Science Grid

### Methods

#### How to submit jobs to the Open Science Grid:

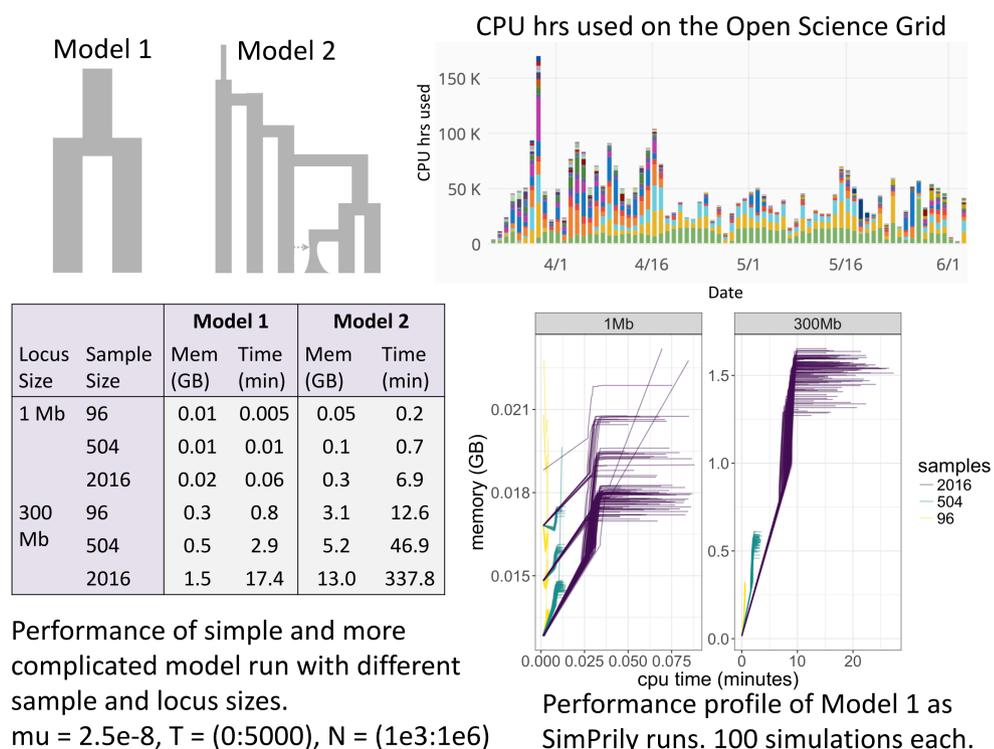
```
./submit param_file.txt model_file.csv array_template genetic_map number_jobs
```

#### High throughput workflow

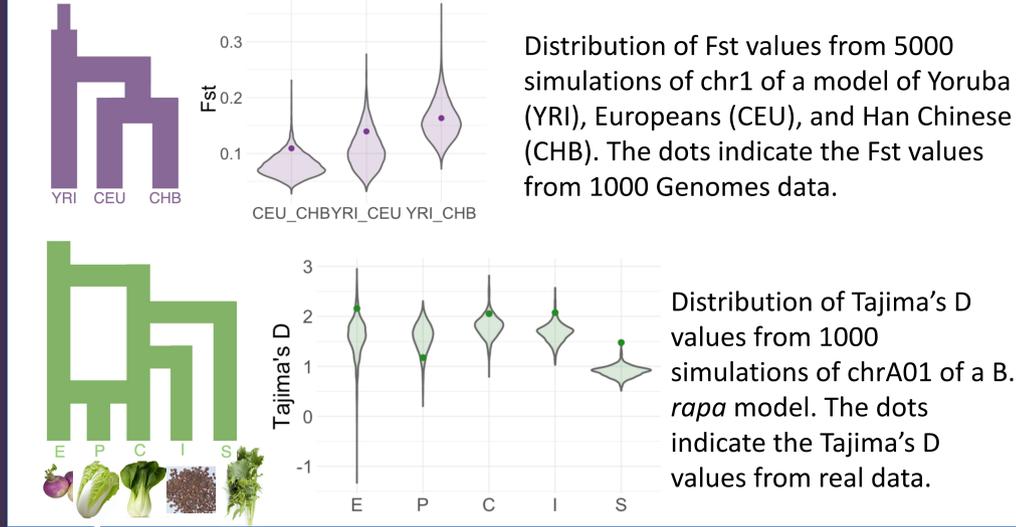


### Results

#### Performance



#### Examples



### Conclusions

#### Major Benefits

- No experience with HTC required - Ready to use HTC workflow
- Minimal storage required – simulations are not printed

#### Future Work

- Options of a variety of simulators
- Serial or parallel multilocus simulations

#### References

- <sup>1</sup>Chen, G. et al. 2009. *Genome research*, 136-142.
- <sup>2</sup>Deelman, E. et al. 2015. *Future Generation Computer Systems*, 46, 17-35.
- <sup>3</sup>Merchant, N. et al. 2016. *Plos Biology*, 14(1), 1-9.
- <sup>4</sup>Pordes, R. et al. 2007. *Journal of Physics: Conference Series*, 78, 012057.
- <sup>5</sup>Qi, X. et al. 2017. *Molecular Ecology*, 26:3373-3388.
- <sup>6</sup>Quinto-Cortés et al. 2017. Submitted
- <sup>7</sup>The 1000 Genomes Project Consortium 2015. *Nature*, 526, 68-74.

#### Acknowledgments

We would like to thank Mats Rynge for his valuable help with setting up the Pegasus workflow and running it on the Open Science Grid.

This material is based upon work supported by the National Science Foundation under Award Numbers DBI-0735191 and DBI-1265383. URL: [www.cyverse.org](http://www.cyverse.org). Pegasus is funded by The National Science Foundation under OAC S12-SSI program, grant \#1664162. This research was done using resources provided by the Open Science Grid, which is supported by the National Science Foundation and the U.S. Department of Energy's Office of Science.