

# **Detecting Selection from Longitudinal Deep Sequencing Data in Viruses**

**Bhavin Khatri**

*bkhatri@imperial.ac.uk*

**in collaboration with Matteo Fumagalli**

*m.fumagalli@imperial.ac.uk*

**Talia Al-Mushadani**

*ta1915@ic.ac.uk*

# 1 Keywords

2 *Evolutionary Virology, Population Genetics, Natural Selection,*  
3 *Machine Learning, Norovirus*

## 4 2 Introduction to Project & Proposed Questions

5 Currently research into how the Norovirus genome evolves via selection  
6 has been limited in scope (Steyer et al., 2018) due to the lack of longitu-  
7 dinal deep-sequencing data on the virus and the lack of a simple methods  
8 to assess where in the genome selection is occurring. This project will be  
9 using longitudinal deep-sequencing data of Norovirus in patients, provided  
10 by Prof. Judith Breuer at UCL/GOSH, and will compare two streamlined  
11 identification methodologies. The two main methods that will be compared  
12 are a heuristic approach (Khatri, 2016) and neural networks, thereby pro-  
13 viding a contrast between methods built on theory vs methods which are  
14 biologically uninformed. As such this project has two main questions that it  
15 wishes to address:  
16 1) Which method is most effective at accurately predicting selection?  
17 2) Which loci are under selection in the Norovirus genome in each individ-  
18 ual?

## 19 3 Proposed Methods

20 Initially the heuristic approach algorithm (Khatri, 2016) will be tested using  
21 simulation longitudinal deep-sequencing data for a single loci with a known  
22 selection strength, employing a multiple hypothesis procedure with false  
23 discovery rates (Benjamini and Hochberg, 1995), when identifying whether  
24 the peaks in nucleotide frequency in the loci over time are characteristic of  
25 selection. A recurrent neural network (RNN) will then be created (using the  
26 same simulated data) to be able to recognise the strength of selection in the  
27 single loci.

28 Both of these methods will then be extended to be able to look at 2 loci.  
29 Finally these methods would be expanded to take into account an entire  
30 section of viral DNA, which would involve producing another neural network  
31 consisting, possibly, of exchangeable convolutional layers (Chan et al., 2018)

32 with Markov intuition, that are capable of using image data for identifying se-  
33 lection.

34 The accuracy of the methods at each stage will be assessed using simu-  
35 lated data, before applying the most accurate methods to the real Norovirus  
36 patient data, to identify which specific loci are under selection.

37 If time permits, an additional numerical path integrating approach may also  
38 be considered, which is based upon theory commonly found in computa-  
39 tional physics (the method is outlined in an as yet unpublished paper by B.  
40 Khatri).

## 41 **4 Anticipated Outcomes**

42 I anticipate that at the end of this project I will have been able to identify loci  
43 in the Norovirus genomes from each patient which are under selection. This  
44 would enable further research into the identification of the specific genes  
45 that are being selected to evolve whilst Norovirus is infecting a host, al-  
46 lowing for the generation of new hypotheses which can be investigated via  
47 molecular assays.

## 48 **5 Project Feasibility & Proposed Timeline**

49 The project is feasible as the heuristic approach has already been shown  
50 to work on single loci simulations (Khatri, 2016), and the theory behind the  
51 algorithm is applicable to multiple loci. Also, most of the machine learning  
52 techniques that will be employed have already been applied to longitudinal  
53 data (Choi et al., 2017; Wang et al., 2018), although some changes to the  
54 current exchangeable neural network structure (Chan et al., 2018) would  
55 need to be made for the purposes of this project. This however is all de-  
56 pendent on the amount of time taken by each step, as to the number of loci  
57 that can be considered at once.

## 58 **6 Itemised Budget**

### 59 **Research-related travel: £420**

60 As I do not live at Silwood, and will be occupying desk space at Silwood

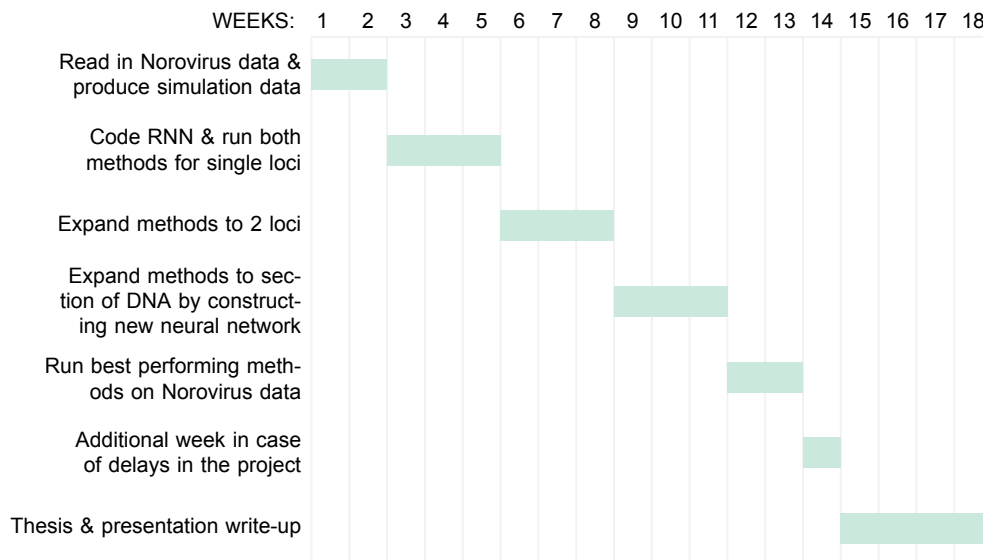


Figure 1: Gantt chart for this project

61 instead of at the Crick Institute, where my main supervisor is mainly based,  
 62 an outlay of £420, which accounts to 35 days return train fare (half the num-  
 63 ber of days that I expect to travel into Silwood over the course of 14 weeks)  
 64 would be of great assistance to me.

## 65 7 Cited References

- 66 Benjamini, Y. and Hochberg, Y. (1995), 'Controlling the False Discovery  
 67 Rate: A Practical and Powerful Approach to Multiple Testing'.  
 68 **URL:** <https://www.jstor.org/stable/2346101>
- 69 Chan, J., Perrone, V., Spence, J. P., Jenkins, P. A., Mathieson, S. and  
 70 Song, Y. S. (2018), A Likelihood-Free Inference Framework for Popula-  
 71 tion Genetic Data using Exchangeable Neural Networks, Technical report.  
 72 **URL:** <https://github.com/popgenmethods/>
- 73 Choi, E., Schuetz, A., Stewart, W. F. and Sun, J. (2017), 'Using recurrent  
 74 neural network models for early detection of heart failure onset.', *Journal*  
 75 *of the American Medical Informatics Association* : *JAMIA* **24**(2), 361–  
 76 370.

- 77 **URL:** <http://www.ncbi.nlm.nih.gov/pubmed/27521897>  
78 <http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=PMC5391725>
- 79 Khatri, B. S. (2016), 'Quantifying evolutionary dynamics from variant-  
80 frequency time series', *Scientific Reports* **6**(1), 32497.  
81 **URL:** <http://www.nature.com/articles/srep32497>
- 82 Steyer, A., Konte, T., Sagadin, M., Kolenc, M., Škoberne, A., Germ, J.,  
83 Dovč-Drnovšek, T., Arnol, M. and Poljšak-Prijatelj, M. (2018), 'Intrahost  
84 Norovirus Evolution in Chronic Infection Over 5 Years of Shedding in a  
85 Kidney Transplant Recipient', *Frontiers in Microbiology* **9**, 371.  
86 **URL:** <http://journal.frontiersin.org/article/10.3389/fmicb.2018.00371/full>
- 87 Wang, T., Qiu, R. G. and Yu, M. (2018), 'Predictive Modeling of the Progres-  
88 sion of Alzheimer's Disease with Recurrent Neural Networks', *Scientific*  
89 *Reports* **8**(1), 9161.  
90 **URL:** <http://www.nature.com/articles/s41598-018-27337-w>