

1

2

## Origin time and epidemic dynamics of the 2019 novel coronavirus

3

4

Chi Zhang<sup>1,2,\*</sup> and Mei Wang

5

January 25, 2020

6

7

<sup>1</sup>Key Laboratory of Vertebrate Evolution and Human Origins, Institute of Vertebrate Paleontology and Paleoanthropology, Chinese Academy of Sciences, Beijing 100044, China

8

9

<sup>2</sup>Center for Excellence in Life and Paleoenvironment, Chinese Academy of Sciences, Beijing 100044, China

10

11

\*Corresponding author: E-mail: zhangchi@ivpp.ac.cn

12

13

### Abstract

14

The 2019 novel coronavirus (2019-nCoV) have emerged from Wuhan, China. Studying the epidemic dynamics is crucial for further surveillance and control of the outbreak. We employed a Bayesian framework to infer the time-calibrated phylogeny and the epidemic dynamics represented by the effective reproductive number ( $R_e$ ) changing over time from the genomic sequences available from GISAID. The origin time is estimated to be December 17, 2019 (95% CI: December 5, 2019 – December 23, 2019). The median estimate of  $R_e$  ranges from 0.2 to 2.2 and changes drastically over time. This study provides an early insight of the 2019-nCoV epidemic.

22

23

### Introduction

24

An outbreak of a novel coronavirus (2019-nCoV) was reported in Wuhan, a city in central China (WHO). Coronaviruses cause diseases range from common cold to severe pneumonia.

25

26

Two fatal coronavirus epidemics over the last two decades were severe acute respiratory

27 syndrome (SARS) in 2003 and Middle East respiratory syndrome (MERS) in 2012 (WHO).  
28 Human to human transmission has been confirmed for this new type of coronavirus (Wang et  
29 al. 2020) and more than 1,000 cases have been reported as of January 25, 2020  
30 (<https://bnonews.com/index.php/2020/01/the-latest-coronavirus-cases>).

31 Studying the virus epidemic dynamics is crucial for further surveillance and control of the  
32 outbreak. Phylogeny of the viruses is a proxy of the transmission chain. Early studies of 2019-  
33 nCoV have focused on their molecular features and phylogenetic relationship with the close  
34 relatives (Zhou et al. 2020; Chan et al. 2020). The phylogenetic analyses typically ignored the  
35 sampling times and measured branch lengths by expected number of substitutions per site.  
36 They are also lacking a stochastic process to model the epidemic dynamics over time. However,  
37 both the timing and dynamics are critical to understand the early outbreak of 2019-nCoV.  
38 Furthermore, various sources of information and uncertainties are hard to be integrated in the  
39 analyses without employing a Bayesian approach.

40 In this study, we used the birth-death skyline serial (BDSS) model (Stadler et al. 2013) to  
41 infer the phylogeny, divergence times and epidemic dynamics of 2019-nCoV. This approach  
42 takes the genomic sequences and sampling times of the viruses as input, and co-estimates the  
43 phylogeny and key epidemic parameters in a Bayesian framework. To our knowledge, this is  
44 the first study to perform such estimation on 2019-nCoV.

## 45 **Results and Discussion**

46 The phylogeny in Figure 1 shows the divergence times and relationships of the 24  
47 BetaCoV viruses. The samples from Wuhan form a paraphyletic group, while the rest of the  
48 samples form a monophyletic clade. These two clades are not divergent from each other as  
49 their sequences are quite similar, which indicates that the outbreak is still in an early stage. The  
50 patients from Guangdong and Zhejiang had traveled from Wuhan and they had been infected  
51 before travelling (Chan et al. 2020). Note that this phylogeny is a maximum clade credibility  
52 (MCC) tree summarized from the posterior samples, which represents a best estimate of the  
53 topology. The probabilities in most clades are lower than 0.5 and would form polytomies if  
54 summarized as a 50% majority-rule consensus tree (GISAID).

55 The time of origin is estimated with median 31.9 days and 95% highest posterior density  
56 (HPD) interval from 26.3 to 43.8 days before the date of the latest sample (January 18, 2020),  
57 that is, December 17, 2019 (December 5, 2019 – December 23, 2019) (Table 1, ten intervals).  
58 This is in agreement with the symptom onset reported by WHO, showing again an early phase  
59 of the epidemic.

60 We investigate the epidemic dynamics of 2019-nCoV by estimating  $R_e$  changing over time  
61 in the BDSS model.  $R_e > 1.0$  means that the number of cases are increasing and the epidemic  
62 is growing, whereas  $R_e < 1.0$  means that the epidemic is declining and will die out. Interestingly,  
63 the epidemic has an early boost with  $R_e$  at around 2.2, then decreases dramatically to 0.2, and  
64 increases again to about 1.7 during the last phase (Table 1). In general, this is in agreement  
65 with some other studies reporting  $R_e$  ranging from 1.4 to 5.5 (Read et al. 2020; Zhao et al. 2020;  
66 Riou and Althaus 2020). In comparison, the estimated  $R_e$  was 2.7 to 3.6 for SARS during the  
67 precontrol phase in Hong Kong (Riley et al. 2003; Wallinga and Teunis 2004). Dividing  $R_e$   
68 into ten intervals rather than three give us a better resolution (Figure 2). This drastic change  
69 could reflect the epidemic to some extent but is likely sensitive to the virus sampling. Keep in  
70 mind that we used only 24 samples in our analysis, which is less than 2% of the reported  
71 number of infected patients, thus one needs to be cautious when interpreting this result. With  
72 more viruses sequenced, we would expect more reliable estimates which would provide better  
73 insights into the epidemic of 2019-nCoV.

74 Overall, this study provides an early insight of the 2019-nCoV epidemic by inferring key  
75 epidemiological parameters from the virus sequences. Such estimates would help public health  
76 officials to coordinate effectively to control the outbreak.

## 77 **Methods**

78 The original data downloaded from GISAID (<https://www.gisaid.org>) consists of 26  
79 sequences of 2019-nCoV (as of January 24, 2020). We excluded two outliers, one is the virus  
80 from Kanagawa, Japan (EPI\_ISL\_402126) which contained only a small segment of 369bp,  
81 another is from Wuhan, China (EPI\_ISL\_403928) which contains suspiciously many  
82 mutations and has genetic distance about ten times longer than the rest of the sequences

83 (GISAID). Sequence alignment was done using MUSCLE (Edgar 2004), resulting in a total  
84 length of 29904bp for the whole genome. The sampling times of the viruses ranged from  
85 December 24, 2019 (EPI\_ISL\_402123) to January 18, 2020 (EPI\_ISL\_403937) and they  
86 were used as fixed ages (in unit of years) in subsequent analysis.

87 We used the BDSS model (Stadler et al. 2013) implemented in the BDSKY package for  
88 BEAST 2 (Bouckaert et al. 2019) to infer the phylogeny, divergence times and epidemic  
89 dynamics of 2019-nCoV. The model has an important epidemiological parameter, the effective  
90 reproductive number  $R_e$ , defined as the number of expected secondary infections caused by an  
91 infected individual during the epidemic. The model allows  $R_e$  to change over time, making it  
92 feasible to estimate its dynamics (Stadler et al. 2013). The BDSS process starts from the origin  
93 time  $t_0$ , which was assigned a lognormal(-2, 1.5) prior with median 0.135 (years before the  
94 latest sampling time). Time from the origin to the latest sample was divided into either three or  
95 ten equally spaced intervals in which  $R_e$  was varied and estimated individually in each interval.  
96 The prior for  $R_e$  was a lognormal(0, 1.25) distribution with median 1.0 and 95% quantiles  
97 between 0.13 and 7.82. The other two parameters are the becoming noninfectious rate  $\delta$  and  
98 sampling proportion  $p$ , which were assumed constant over time.  $\delta$  was given a lognormal (2,  
99 1.25) prior with median 7.39 and mean 16.1, expecting the infectious period of an individual  
100 ( $1/\delta$ ) to be about a month. The sampling proportion of infected individuals  $p$  was a beta(1, 9)  
101 distribution with mean 0.1.

102 We used the lognormal independent relaxed clock (Drummond et al. 2006; Rannala and  
103 Yang 2007) to model evolutionary rate variation along the branches. The mean clock rate  $r$  was  
104 assigned a gamma(2, 0.0005) prior with mean of 0.001 substitutions per site per year and the  
105 standard deviation  $\sigma$  was a gamma (0.54, 0.38) prior with mean 0.2 and median 0.1 by default.  
106 The substitution model used was HKY+ $\Gamma_4$  (Hasegawa et al. 1985; Yang 1994) in which the  
107 transition-transversion rate ratio  $\kappa$  was set a lognormal(1, 1.25) prior and the gamma shape  
108 parameter  $\alpha$  was an exponential(1) prior.

109 The analysis was performed in the BEAST 2 platform (Bouckaert et al. 2019). We ran 150  
110 million MCMC iterations and sampled every 5000 iterations. The first 30% samples were  
111 discarded as burn-in. Convergence was diagnosed in Tracer (Rambaut et al. 2018) to confirm  
112 that independent runs gave consensus results and all parameters had effective sample size (ESS)

113 larger than 100. The remaining 70% samples were used to build the maximum clade credibility  
114 (MCC) tree and to summarize the parameter estimates. The common ancestor heights were  
115 used to annotate the clade ages in the MCC tree.

116

## 117 **Acknowledgments**

118 CZ is supported by the 100 Young Talents Program of Chinese Academy of Sciences.

119

## 120 **References**

121 Bouckaert R., Vaughan T.G., Barido-Sottani J., Duchêne S., Fourment M., Gavryushkina A.,  
122 Heled J., Jones G., Kühnert D., De Maio N., Matschiner M., Mendes F.K., Müller N.F.,  
123 Ogilvie H.A., Plessis du L., Poppinga A., Rambaut A., Rasmussen D., Siveroni I., Suchard  
124 M.A., Wu C.-H., Xie D., Zhang C., Stadler T., Drummond A.J. 2019. BEAST 2.5: An  
125 advanced software platform for Bayesian evolutionary analysis. *PLoS Comput. Biol.*  
126 15:e1006650.

127 Chan J.F.-W., Yuan S., Kok K.-H., To K.K.-W., Chu H., Yang J., Xing F., Liu J., Yip C.C.-  
128 Y., Poon R.W.-S., Tsoi H.-W., Lo S.K.-F., Chan K.-H., Poon V.K.-M., Chan W.-M., Ip  
129 J.D., Cai J.-P., Cheng V.C.-C., Chen H., Hui C.K.-M., Yuen K.-Y. 2020. A familial  
130 cluster of pneumonia associated with the 2019 novel coronavirus indicating person-to-  
131 person transmission: a study of a family cluster. *The Lancet*.

132 Drummond A.J., Ho S.Y.W., Phillips M.J., Rambaut A. 2006. Relaxed phylogenetics and  
133 dating with confidence. *PLoS Biol.* 4:e88.

134 Edgar R.C. 2004. MUSCLE: multiple sequence alignment with high accuracy and high  
135 throughput. *Nucleic Acids Res.* 32:1792–1797.

136 GISAID. 2020 Coronavirus. <https://www.gisaid.org/CoV2020>.

137 Hasegawa M., Kishino H., Yano T. 1985. Dating of the human-ape splitting by a molecular  
138 clock of mitochondrial DNA. *J. Mol. Evol.* 22:160–174.

139 Rambaut A., Drummond A.J., Xie D., Baele G., Suchard M.A. 2018. Posterior  
140 summarization in Bayesian phylogenetics using Tracer 1.7. *Syst. Biol.* 67:901–904.

141 Rannala B., Yang Z. 2007. Inferring speciation times under an episodic molecular clock.  
142 *Syst. Biol.* 56:453–466.

- 143 Stadler T., Kühnert D., Bonhoeffer S., Drummond A.J. 2013. Birth-death skyline plot reveals  
144 temporal changes of epidemic spread in HIV and hepatitis C virus (HCV). *Proc. Natl.*  
145 *Acad. Sci. USA.* 110:228–233.
- 146 Wang C., Horby P.W., Hayden F.G., Gao G.F. 2020. A novel coronavirus outbreak of global  
147 health concern. *The Lancet.*
- 148 World Health Organization (WHO). Novel Coronavirus – China.  
149 <https://www.who.int/csr/don/12-january-2020-novel-coronavirus-china>. (date accessed:  
150 January 25, 2020)
- 151 World Health Organization (WHO). Coronavirus. [https://www.who.int/health-](https://www.who.int/health-topics/coronavirus)  
152 [topics/coronavirus](https://www.who.int/health-topics/coronavirus). (date accessed: January 25, 2020)
- 153 Yang Z. 1994. Maximum likelihood phylogenetic estimation from DNA sequences with  
154 variable rates over sites: approximate methods. *J. Mol. Evol.* 39:306–314.
- 155 Zhou P., Yang X.-L., Wang X.-G., Hu B., Zhang L., Zhang W., Si H.-R., Zhu Y., Li B.,  
156 Huang C.-L., Chen H.-D., Chen J., Luo Y., Guo H., Jiang R.-D., Liu M.-Q., Chen Y.,  
157 Shen X.-R., Wang X., Zheng X.-S., Zhao K., Chen Q.-J., Deng F., Liu L.-L., Yan B.,  
158 Zhan F.-X., Wang Y.-Y., Xiao G., Shi Z.-L. 2020. Discovery of a novel coronavirus  
159 associated with the recent pneumonia outbreak in humans and its potential bat origin.  
160 bioRxiv.
- 161
- 162

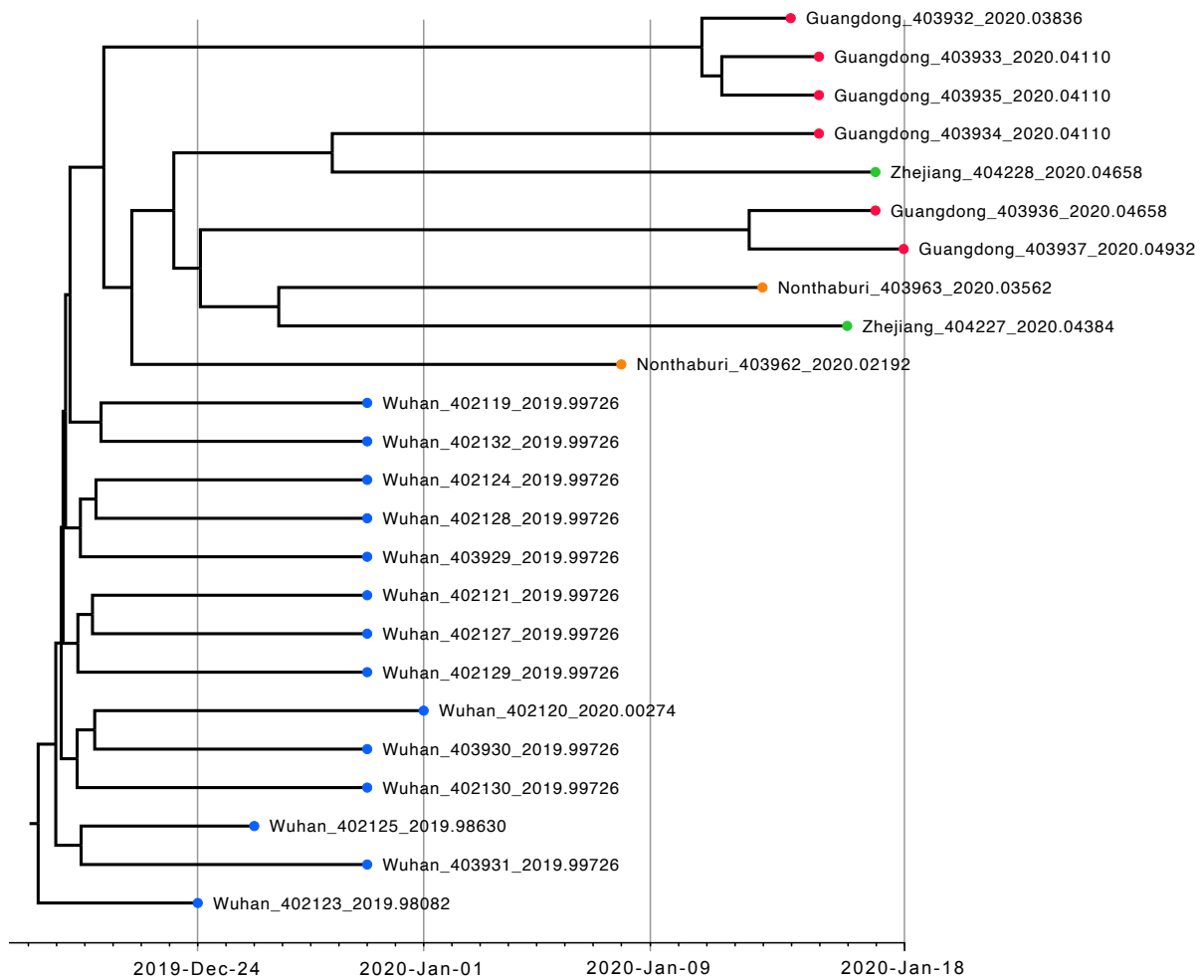
163 Table 1. Posterior estimates (median and 95% HPD interval) of key model parameters

	three time intervals	ten time intervals
$t_0$	0.0812 (0.076, 0.099)	0.0873 (0.072, 0.120)
$R_e$	2.18 (1.31, 3.64)	1.32 (0.013, 5.92)
	0.21 (0.014, 0.55)	1.72 (0.015, 5.58)
	1.71 (1.11, 2.46)	1.07 (0.015, 4.50)
		1.79 (0.015, 3.51)
		0.28 (0.006, 2.61)
		0.31 (0.006, 1.47)
		0.59 (0.014, 1.93)
	0.59 (0.012, 2.10)	
	2.54 (0.52, 4.87)	
	0.84 (0.014, 2.73)	
$\delta$	148.31 (32.44, 293.12)	183.25 (15.45, 451.97)
$p$	0.038 (0.00052, 0.20)	0.0099 (8.8E-7, 0.18)
$r$	0.0019 (0.00095, 0.0034)	0.0015 (0.00067, 0.0030)
$\sigma$	0.85 (1.2E-7, 1.44)	0.79 (2.3E-7, 1.40)
$\kappa$	7.41 (2.68, 15.91)	7.38 (2.74, 16.05)
$\alpha$	0.68 (0.0011, 2.93)	0.73 (0.0015, 3.00)

164 Note: time unit is years.

165

166



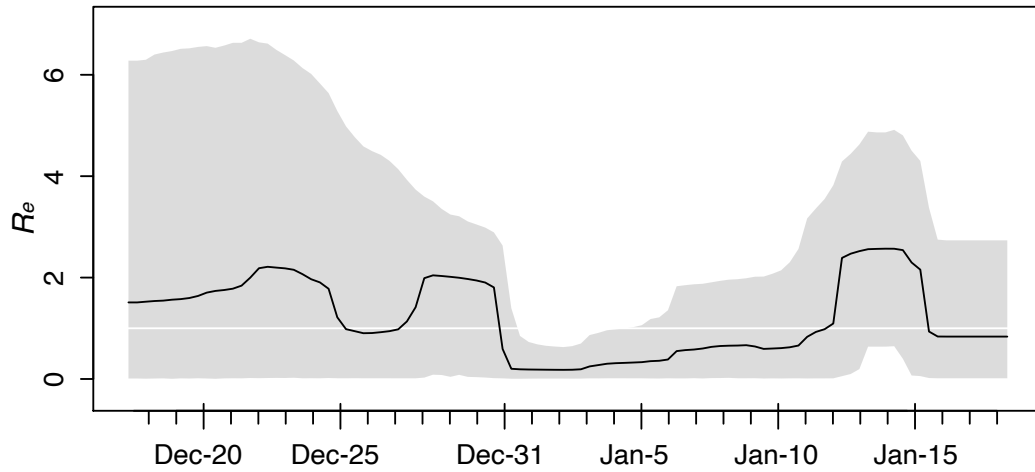
167

168 Figure 1. Maximum clade credibility (MCC) tree summarized from the MCMC sample. The

169 common ancestor heights were used to annotate the clade ages.

170





171

172 Figure 2. Skyline plot of the effective reproductive number ( $R_e$ ) changing over time. The solid

173 line indicates the median and the gray area indicates the 95% HPD interval.

174

175